# 4 Appendix

#### 4.1 Additional related works

**Machine unlearning and memorization.** Many unlearning methods are proposed to effectively erase information of selected samples. Several basic but well-known methods such as random labeling of forget set (Graves et al., 2020) and explicit gradient ascent on the forget set (Warnecke et al., 2023) lay foundation for current unlearning methods. More recent works extend on those works to improve overall performance of unlearning. For example, SCRUB (Kurmanji et al., 2023) simultaneously perform gradient ascent on forget set and gradient descent in the retain set to better preserve performance of retain during unlearning. Influence based (Izzo et al., 2021) unlearning propose idea that takes into account of the Hessian information of datasets to perform update of the model weights. Saliency Unlearning (Fan et al., 2024) identity weights that react strongly to forget set through magnitude of gradient and perform unlearning only on those weight to achieve better performance. There are several theoretical studies about unlearning through the lens of the differential privacy and provide performance guarantee. For example, Langevin Unlearning Chien et al. (2025) study unlearning with privacy guarantee through projected noisy gradient descent. Sekhari et al. (2021) studies unlearning problem and provide performance guarantee and the corresponding sample complexity. There are also works discussing relationship between memorization and generalization. Attias et al. (2024) discuss the fundamental trade-off between generalization and memorization under information theory framework. Carlini et al. (2019) discuss different metrics for identifying sample of different type (memorized, prototypical and so on). Feldman (2021) provide theoretical and experimental analysis saying the memorization is necessary to achieve optimal performance. There are also several works studying memorization with different tasks and model architectures (Biderman et al. (2023); Li et al. (2025); Prashanth et al. (2025)).

### 4.2 Lemmas and proofs

Definition 6 (Full forget Hessian and retain Hessian).

$$H_R = \frac{1}{n_r} \sum_{r \in D_r} H_r, \ H_F = \frac{1}{n_r} \sum_{f \in D_f} H_f,$$
 (13)

**Lemma 4.1.**  $l_1 - l_2$  norm inequality: For any  $x \in \mathbb{R}$ ,  $||x||_2 \le ||x||_1 \le \sqrt{d}||x||_2$ 

**Lemma 4.2.** Binomial coefficient: For all  $n, k \in \mathbb{N}$  such that  $k \leq n$ , the binomial coefficients satisfy that

$$\binom{n}{k} = \binom{n-1}{k-1} + \binom{n-1}{k}.\tag{14}$$

**Lemma 4.3.** For any matrix  $M \in \mathbb{R}^{n \times n}$ ,  $||M||_F \le ||M||_{S_1} \le \sqrt{n}||M||_F$ , where  $||M||_{S_p}$  is p norm of the spectrum of M, and the inequality is obtain through  $l_1 - l_2$  norm inequality.

**Lemma 4.4.** For matrices  $M_1...M_k \in \mathbb{R}^{n \times n}$ ,  $\text{Tr}[M_1...M_k] \leq ||M_1...M_k||_{S_1}$  (see Bhatia (2013))

**Lemma 4.5** (Vershynin (2018)). Consider n random gaussian vectors  $x_1...x_n$  sampled i.i.d from  $N(0, \sigma^2 I_d)$ , there exist a constant  $C_1$  such that with probability  $1 - \delta$ ,

$$\sum_{i=1}^{n} \|x_i\| \le n\sqrt{2}\sigma \frac{\Gamma((d+1)/2)}{\Gamma(d/2)} + \sqrt{\frac{n\sigma^2}{C_1}\log(\frac{2}{\delta})} \quad \textit{for } n, d \textit{ large enough}. \tag{15}$$

**Lemma 4.6** (Vershynin (2018)). Consider n random gaussian vectors  $x_1...x_n$  sampled i.i.d from  $N(0, \sigma^2 I_d)$ , there exist a constant  $C_1$  such that with probability  $1 - \delta$ ,

$$\sum_{i=1}^{n} \|x_i\|^2 \le n\sigma^2 d + \sqrt{\frac{n\sigma^4 d}{C_2} \log(\frac{2}{\delta})} \quad \text{for } n, d \text{ large enough.}$$
 (16)

### 4.3 Proof of Lemma 2.1

**Lemma 4.7** (Restated). Consider the unlearning update operator  $J_k$  defined in (3). Define a sequence of PSD matrices  $\{N_k\}_{k\geq 0}$  by  $N_0=I$  and for  $k\geq 1$ :

$$N_k = C_f \sum_{i \in D_f} H_i N_{k-1} H_i + C_r \sum_{i \in D_r} H_i N_{k-1} H_i,$$
(17)

with  $C_r$ ,  $C_f$  as given in Definition 2. Also let  $M_k = J^{2k} + N_k$ . Then:

- 1. (Lower bound)  $\mathbb{E}\operatorname{Tr}(J_0^T\cdots J_{k-1}^TJ_{k-1}\cdots J_0) \geq \operatorname{Tr}(M_k)$ . Moreover, if  $\operatorname{Tr}(N_k)\to\infty$  as  $k\to\infty$ , then  $\mathbb{E}\|w_k\|^2\to\infty$  as well.
- 2. (Upper bound) If at each step  $J_k$  is spectrally bounded as  $(1 \epsilon)I \succeq J \succeq -(1 \epsilon)I$  for some  $\epsilon \in (0, 1)$  (i.e. all eigenvalues of J lie in  $[-(1 \epsilon), 1 \epsilon]$ ), then

$$\mathbb{E} \operatorname{Tr} (J_0^T \cdots J_{k-1}^T J_{k-1} \cdots J_0) \leq \sum_{r=0}^{k-1} {k \choose r} (1-\epsilon)^{2(k-r)} \operatorname{Tr} (N_r).$$

If in addition  $\operatorname{Tr}(N_r) \leq \epsilon$  for all r, then  $\mathbb{E}||w_k||^2 \to 0$  as  $k \to \infty$  (the unlearning update converges in mean square).

*Proof.* As we are taking the expectation value over the calculation, we can effectively transform the  $J_k$  into following with random variables involved:

$$J_{k} = (I - \eta(1 - \alpha)\frac{1}{B}\sum_{i \in D_{r,k}} H_{i} + \eta\alpha\frac{1}{B}\sum_{i \in D_{f,k}} H_{i}) = (I - \eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_{r}} x_{r}H_{r} + \eta\alpha\frac{1}{B}\sum_{i \in D_{f}} x_{f}H_{f}),$$
(18)

where  $x_r, x_f$  are the corresponding Bernoulli random variables with probability  $P(x_r=1)=\frac{B}{n_r}$  and  $P(x_f=1)=\frac{B}{n_f}$  and 0 otherwise.

To initiate the first step in characterize the difference between the unlearning and usually learning process, we first calculate the  $E[J_1^T J_1]$  as follows:

$$E[J_{1}^{T}J_{1}] = E[(I - \eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_{r}}x_{r}H_{r} + \eta\alpha\frac{1}{B}\sum_{i \in D_{f}}x_{f}H_{f})^{T}(I - \eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_{r}}x_{r}H_{r} + \eta\alpha\frac{1}{B}\sum_{i \in D_{f}}x_{f}H_{f})]$$

$$= E[(I - \eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_{r}}x_{r}H_{r})^{T}(I - \eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_{r}}x_{r}H_{r})$$

$$+2(I - \eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_{r}}x_{r}H_{r})^{T}(\eta\alpha\frac{1}{B}\sum_{i \in D_{f}}x_{f}H_{f}) + (\eta\alpha\frac{1}{B}\sum_{i \in D_{f}}x_{f}H_{f})^{T}(\eta\alpha\frac{1}{B}\sum_{i \in D_{f}}x_{f}H_{f})].$$
(19)

Here, we separate the above equation into three part and take the expectation accordingly:

$$E[(I - \eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_{r}} x_{r}H_{r})^{T}(I - \eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_{r}} x_{r}H_{r})]$$

$$= E[(I - 2\eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_{r}} x_{r}H_{r} + \eta^{2}(1 - \alpha)^{2}\frac{1}{B}^{2}\sum_{r \in D_{r}} x_{r}H_{r}\sum_{r \in D_{r}} x_{r}H_{r})],$$

$$= I - 2\eta(1 - \alpha)\frac{1}{n_{r}}\sum_{r \in D_{r}} H_{r} + E[\eta^{2}(1 - \alpha)^{2}(\frac{1}{B})^{2}\sum_{r \in D_{r}} x_{r}H_{r}\sum_{r \in D_{r}} x_{r}H_{r}],$$

$$= I - 2\eta(1 - \alpha)\frac{1}{n_{r}}\sum_{r \in D_{r}} H_{r} + E[\eta^{2}(1 - \alpha)^{2}(\frac{1}{B})^{2}\sum_{r' \in D_{r}} \sum_{r' \in D_{r}} x_{r'}x_{r}H_{r'}H_{r}],$$

$$= I - 2\eta(1 - \alpha)\frac{1}{n_{r}}\sum_{r \in D_{r}} H_{r} + \eta^{2}(1 - \alpha)^{2}(\frac{1}{n_{r}})^{2}(\sum_{r \in D_{r}} H_{r})^{2} + \eta^{2}(1 - \alpha)^{2}\frac{1}{n_{r}}(\frac{1}{B} - \frac{1}{n_{r}})\sum_{r \in D_{r}} H_{r}^{2},$$

$$= I - 2\eta(1 - \alpha)H_{R} + \eta^{2}(1 - \alpha)^{2}H_{R}^{2} + \eta^{2}(1 - \alpha)^{2}\frac{1}{n_{r}}(\frac{1}{B} - \frac{1}{n_{r}})\sum_{r \in D_{r}} H_{r}^{2},$$

$$= (I - \eta(1 - \alpha)H_{R})^{2} + \eta^{2}(1 - \alpha)^{2}\frac{1}{n_{r}}(\frac{1}{B} - \frac{1}{n_{r}})\sum_{r \in D_{r}} H_{r}^{2}.$$

$$(20)$$

The random variables  $x_r$  are independent to each other but not itself and therefore there exist one additional terms in the final line. Also, Compared to the original sgd there exists additional multiplication of the  $(1-\alpha)^2$ . Next, we move on to the interaction term:

$$E[2(I - \eta(1 - \alpha)\frac{1}{B}\sum_{r \in D_r} x_r H_r)^T (\eta \alpha \frac{1}{B}\sum_{i \in D_f} x_f H_f)] = 2(I - \eta(1 - \alpha)H_R)^T (\eta \alpha H_F).$$
 (21)

We can directly formulate this as above due to the fact that we assume the sampling process of retain set and forget set to be independent. Last, the term arising due to the forget set:

$$E[(\eta \alpha \frac{1}{B} \sum_{f \in D_f} x_f H_f)^T (\eta \alpha \frac{1}{B} \sum_{f \in D_f} x_f H_f)] = \eta^2 \alpha^2 H_F^2 + \eta^2 \alpha^2 \frac{1}{n_f} (\frac{1}{B} - \frac{1}{n_f}) \sum_{f \in D_f}^{n_f} H_f^2. \quad (22)$$

We then integrate the three part and reformulate the Jacobian:

$$E[J_{1}^{T}J_{1}] = (I - \eta(1 - \alpha)H_{R})(I - \eta(1 - \alpha)H_{R}) + 2(I - \eta(1 - \alpha)H_{R})^{T}(\eta\alpha H_{F}) + \eta^{2}\alpha^{2}H_{F}^{2}$$

$$+ \eta^{2}\alpha^{2}\frac{1}{n_{f}}(\frac{1}{B} - \frac{1}{n_{f}})\sum_{f \in D_{f}}^{n_{f}}H_{f}^{2} + \eta^{2}(1 - \alpha)^{2}\frac{1}{n_{r}}(\frac{1}{B} - \frac{1}{n_{r}})\sum_{r \in D_{r}}^{n_{r}}H_{r}^{2},$$

$$= (I - \eta(1 - \alpha)H_{R} + \eta\alpha H_{F})(I - \eta(1 - \alpha)H_{R} + \eta\alpha H_{F})$$

$$+ \eta^{2}\alpha^{2}\frac{1}{n_{f}}(\frac{1}{B} - \frac{1}{n_{f}})\sum_{f \in D_{f}}^{n_{f}}H_{f}^{2} + \eta^{2}(1 - \alpha)^{2}\frac{1}{n_{r}}(\frac{1}{B} - \frac{1}{n_{r}})\sum_{r \in D_{r}}^{n_{r}}H_{r}^{2},$$

$$= J^{2} + \eta^{2}\alpha^{2}\frac{1}{n_{f}}(\frac{1}{B} - \frac{1}{n_{f}})\sum_{f \in D_{f}}^{n_{f}}H_{f}^{2} + \eta^{2}(1 - \alpha)^{2}\frac{1}{n_{r}}(\frac{1}{B} - \frac{1}{n_{r}})\sum_{r \in D_{r}}^{n_{r}}H_{r}^{2},$$

$$= J^{2} + \eta^{2}\alpha^{2}\frac{1}{n_{f}}(\frac{1}{B} - \frac{1}{n_{f}})\sum_{f \in D_{f}}^{n_{f}}H_{f}^{2} + \eta^{2}(1 - \alpha)^{2}\frac{1}{n_{r}}(\frac{1}{B} - \frac{1}{n_{r}})\sum_{r \in D_{r}}^{n_{r}}H_{r}^{2},$$

$$(23)$$

where we define  $J = I - \eta(1 - \alpha)H_R + \eta\alpha H_F$ . During the whole work, we will be analyzing on these terms to characterize the behavior of unlearning process.

We use inductive proof for both first and second part of the theory and we begin the proof as following:

### First part:

Base case: k = 1

$$M_1 = J^2 + C_f \sum_{f \in D_f}^{n_f} H_f^2 + C_r \sum_{r \in D_r}^{n_r} H_r^2 = J^2 + N_1 \le E[J_1^T J_1], \tag{24}$$

where the left term match the equation 23 and therefore the basis case is set. Now we go further with inductive step.

## **Inductive step: k-1**

$$\begin{split} E[J_{k}^{T}J_{k-1}^{T}...J_{1}^{T}J_{1}...J_{k-1}] \succeq E[J_{k}^{T}M_{k-1}J_{k}], \\ &= E[(I - \eta(1 - \alpha)\frac{1}{B}\sum_{i \in D_{r}}x_{r}H_{r} + \eta\alpha\frac{1}{B}\sum_{f \in D_{f}}x_{f}H_{f})M_{k-1}(I - \eta(1 - \alpha)\frac{1}{B}\sum_{i \in D_{r}}x_{r}H_{r} + \eta\alpha\frac{1}{B}\sum_{i \in D_{f}}x_{f}H_{f})], \\ &= JM_{k-1}J + C_{f}\sum_{f \in D_{f}}H_{f}M_{k-1}H_{f} + C_{r}\sum_{r \in D_{r}}H_{r}M_{k-1}H_{r}, \\ &= J(J^{2(k-1)} + N_{k-1})J + C_{f}\sum_{f \in D_{f}}H_{f}(J^{2(k-1)} + N_{k-1})H_{f} + C_{r}\sum_{r \in D_{r}}H_{r}(J^{2(k-1)} + N_{k-1})H_{r}, \\ &= J^{2k} + C_{f}\sum_{f \in D_{f}}H_{f}N_{k-1}H_{f} + C_{r}\sum_{r \in D_{r}}H_{r}N_{k-1}H_{r} + JN_{k-1}J + C_{f}\sum_{f \in D_{f}}H_{f}J^{2(k-1)}H_{f} + C_{r}\sum_{r \in D_{r}}H_{r}J^{2(k-1)}H_{r}, \\ &= M_{k} + JN_{k-1}J + C_{f}\sum_{f \in D_{f}}H_{f}J^{2(k-1)}H_{f} + C_{r}\sum_{r \in D_{r}}H_{r}J^{2(k-1)}H_{r}, \\ &\succeq M_{k}. \end{split}$$

The last equality is due to the later three terms are both PSD by assumption as they are symmetric in terms of left and right half of whole multiplication. As we can lower bound through  $M_k$ , diverging of  $N_k$  will lead to  $M_k$  and cause the whole product to diverge.

## Second part:

Base step: k=1.

$$E[J_1^T J_1] = J^2 + N_1 \le (1 - \epsilon)^2 I + N_1 = \sum_{r=0}^{1} {1 \choose r} (1 - \epsilon)^{2(1-r)} N_r.$$
 (26)

The  $J^2$  is bounded by  $(1-\epsilon)^2I$  due to our assumption. Now, we start with the inductive step

**Inductive step: k-1.** 

$$E[J_k^T J_{k-1}^T ... J_1^T J_1 ... J_{k-1}] \leq E[J_k^T \left(\frac{k-1}{r}\right) (1-\epsilon)^{2(k-1-r)} N_r) J_k],$$

$$= J(\sum_{r=0}^{k-1} {k-1 \choose r} (1-\epsilon)^{2(k-1-r)} N_r) J$$

$$+ C_f \sum_{i \in D_f}^{n_f} H_i \left(\sum_{r=0}^{k-1} {k-1 \choose r} (1-\epsilon)^{2(k-1-r)} N_r) H_i + C_r \sum_{i \in D_r}^{n_r} H_i \left(\sum_{r=0}^{k-1} {k-1 \choose r} (1-\epsilon)^{2(k-1-r)} N_r) H_i,$$

$$= J(\sum_{r=0}^{k-1} {k-1 \choose r} (1-\epsilon)^{2(k-1-r)} N_r) J + \sum_{r=0}^{k-1} {k-1 \choose r} (1-\epsilon)^{2(k-1-r)} (C_f \sum_{i \in D_f}^{n_f} H_i N_r H_i + C_r \sum_{i \in D_r}^{n_r} H_i N_r H_i),$$

$$\leq \sum_{r=0}^{k-1} {k-1 \choose r} (1-\epsilon)^{2(k-r)} N_r + \sum_{r=0}^{k-1} {k-1 \choose r} (1-\epsilon)^{2(k-1-r)} N_{r+1},$$

$$= \sum_{r=0}^{k-1} {k-1 \choose r} (1-\epsilon)^{2(k-r)} N_r + \sum_{r=0}^{k-1} {k-1 \choose r} (1-\epsilon)^{2(k-1-r)} N_{r+1},$$

$$= (1-\epsilon)^2 N_o + \sum_{r=1}^{k-1} {k \choose r} N_r + N_k,$$

$$= (1-\epsilon)^2 N_o + \sum_{r=1}^{k-1} {k \choose r} N_r + N_k,$$

$$= \sum_{r=0}^{k} {k \choose r} (1-\epsilon)^{2(k-r)} N_r.$$

$$(27)$$

The first and second inequality is due to the assumption in induction on previous step and we merge the coefficient in the last step through lemma 4.2. Finally, if we further have that  $\text{Tr}[N_r] \leq \epsilon \ \forall r$ , then

$$E[\operatorname{Tr}[J_k^T J_{k-1}^T ... J_1^T J_1 ... J_{k-1} J_{k-1}]],$$

$$= \sum_{r=0}^k \binom{k}{r} (1-\epsilon)^{2k-r} \operatorname{Tr}[N_r],$$

$$\leq \sum_{r=0}^k \binom{k}{r} (1-\epsilon)^{2(k-r)} \epsilon^r,$$

$$\leq ((1-\epsilon)^2 + \epsilon)^k \leq (1-\epsilon)^k,$$
(28)

which will converge to zero when  $k \to \infty$ .

### 4.4 Proof of theorem 2.2

**Theorem 4.8** (Restated). Under the setup of Lemma 2.1, the unlearning process will diverge if the mix-Hessian eigenvalue exceeds a threshold determined by the coherence. In particular, if

$$\lambda_{\max}(D) \geq \frac{\sqrt{2}\sigma}{\eta\left((1-\alpha)n_f\sqrt{\frac{n_r}{B}-1} + \alpha n_r\sqrt{\frac{n_f}{B}-1}\right)},\tag{29}$$

then  $\lim_{k\to\infty} \mathbb{E}||w_k||^2 = \infty$ . Equivalently, condition (7) guarantees the unlearning algorithm will escape the original minima (diverge) due to the stochastic dynamics.

*Proof.* To simplify the notation, we use the following

$$L_k \in \{r, f\}^k$$
 (a string of length  $k$  over the alphabet  $\{r, f\}$ ),  $L_k[i] \mapsto \text{the } i\text{-th symbol of } L_k, \ 1 \le i \le k.$  (30)

We know that based on part one lemma 2.1, we can lower bound the  $N_k$  to lower bound the  $E[\text{Tr}[J_k^T J_{k-1}^T ... J_1^T J_1 ... J_{k-1} J_k]]$ . First, We can write the overall sum as follows:

$$\begin{split} &\operatorname{Tr}[N_{k}] = \sum_{L_{k} \in \{r, f\}^{k}} C_{r}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = f\}} \left( \sum_{a_{k} \in D_{L_{k}[k]}} \sum_{a_{k-1} \in D_{L_{k}[k-1]}} \dots \sum_{a_{1} \in D_{L_{k}[k]}} \right) \operatorname{Tr}[H_{a_{k}} \dots H_{a_{1}} H_{a_{1}} \dots H_{a_{k}}], \\ &= \sum_{L_{k} \in \{r, f\}^{k}} C_{r}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = f\}} \left( \sum_{a_{k} \in D_{L_{k}[k]}} \sum_{a_{k-1} \in D_{L_{k}[k-1]}} \dots \sum_{a_{1} \in D_{L_{k}[k]}} \right) \|H_{a_{k}} \dots H_{a_{1}}\|_{F}^{2}, \\ &\geq \frac{1}{d} \sum_{L_{k} \in \{r, f\}^{k}} C_{r}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = f\}} \left( \sum_{a_{k} \in D_{L_{k}[k]}} \sum_{a_{k-1} \in D_{L_{k}[k-1]}} \dots \sum_{a_{1} \in D_{L_{k}[k]}} \right) \|H_{a_{k}} \dots H_{a_{1}}\|_{S_{1}}^{2}, \\ &\geq \frac{1}{d} \sum_{L_{k} \in \{r, f\}^{k}} C_{r}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = f\}} \left( \sum_{a_{k} \in D_{L_{k}[k]}} \sum_{a_{k-1} \in D_{L_{k}[k-1]}} \dots \sum_{a_{1} \in D_{L_{k}[k]}} \right) \operatorname{Tr}[H_{a_{k}} \dots H_{a_{1}}]^{2}, \\ &\geq \frac{1}{n_{r} n_{f}} \frac{1}{d} \sum_{L_{k} \in \{r, f\}^{k}} C_{r}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} \operatorname{Tr}[H_{a_{L_{k}[k]}} \dots H_{a_{L_{k}[k]}} \dots H_{a_{L_{k}[k]}}]^{2}, \\ &\geq \frac{1}{n_{r} n_{f}} \frac{1}{d} \sum_{a_{r} \in D_{r}} \sum_{a_{f} \in D_{f}} \sum_{1} \frac{1}{2^{k}} \sum_{L_{k} \in \{r, f\}^{k}} \operatorname{Tr}[C_{r}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = f\}} \operatorname{Tr}[H_{a_{L_{k}[k]}} \dots H_{a_{L_{k}[k]}} \dots H_{a_{L_{k}[k]}}]^{2}, \\ &= \frac{1}{n_{r} n_{f}} \frac{1}{d} \sum_{a_{r} \in D_{r}} \sum_{a_{f} \in D_{f}} \sum_{1} \frac{1}{2^{k}} (\operatorname{Tr}[C_{r}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = f\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = f\}} H_{a_{L_{k}[k]}} \dots H_{a_{L_{k}[k]}} \dots H_{a_{L_{k}[k]}}]^{2}, \\ &= \frac{1}{n_{r} n_{f}} \frac{1}{d} \sum_{a_{r} \in D_{r}} \sum_{a_{f} \in D_{f}} \sum_{1} \frac{1}{2^{k}} (\operatorname{Tr}[C_{r}^{\sum_{i=1}^{k} 1\{L_{k}[i] = r\}} C_{f}^{\sum_{i=1}^{k} 1\{L_{k}[i] = f\}} H_{a_{r}} + \frac{C_{f}^{\frac{1}{2}}} H_{a_{r}} + \frac{C_{f}^{\frac{1}{2}}} H_{a_{f}})^{k}}{C_{r}^{\frac{1}{2}} + C_{f}^{\frac{1}{2}}} H_{a_{f}} + \frac{C_{f}^{\frac{1}{2}}} H_{a_{f}} + \frac{C_{f}^{\frac{1$$

For the first and second inequality, we use lemma 4.3 and 4.4. For the third inequality, we reduce the summation to  $\sum_{a_r \in D_r} \sum_{a_f \in D_f}$ . As there are terms without  $D_r$  or  $D_f$  involved, we divided the whole equation by  $n_f n_r$  to ensure inequality. For the forth inequality, we use the lemma 4.1.

Before we try to connect the relationship between the quantity to the above, we first reindex the following:

$$\frac{1}{n_r n_f} \sum_{a_r \in D_r, a_f \in D_f} \frac{C_r^{\frac{1}{2}}}{C_r^{\frac{1}{2}} + C_f^{\frac{1}{2}}} H_{a_r} + \frac{C_f^{\frac{1}{2}}}{C_r^{\frac{1}{2}} + C_f^{\frac{1}{2}}} H_{a_f} = \frac{1}{n_r n_f} \sum_{rf} D_{rf} = D,$$
 (32)

where  $D_{rf} = \frac{C_r^{\frac{1}{2}}}{C_r^{\frac{1}{2}} + C_f^{\frac{1}{2}}} H_r + \frac{C_f^{\frac{1}{2}}}{C_r^{\frac{1}{2}} + C_f^{\frac{1}{2}}} H_f$  and the subscript indicates that summing over corresponding subset (retain and forget set). Now, we proceed to relate different quantities

$$\operatorname{Tr}\left[\left(\frac{1}{n_{r}n_{f}}\sum_{a_{r}\in D_{r}, a_{f}\in D_{f}}\frac{C_{r}^{\frac{1}{2}}}{C_{r}^{\frac{1}{2}}+C_{f}^{\frac{1}{2}}}H_{a_{r}}+\frac{C_{f}^{\frac{1}{2}}}{C_{r}^{\frac{1}{2}}+C_{f}^{\frac{1}{2}}}H_{a_{f}}\right)^{k}\right] = \operatorname{Tr}\left[\left(\frac{1}{n_{r}n_{f}}\right)^{k}\left(\sum_{rf}D_{rf}\right)^{k}\right],$$

$$= \operatorname{Tr}\left[\left(\frac{1}{n_{r}n_{f}}\right)^{k}\sum_{rf_{1}}\sum_{rf_{2}}...\sum_{rf_{k}}D_{rf_{1}}D_{rf_{2}}...D_{rf_{k-1}}D_{rf_{k}}\right],$$

$$\leq d\left(\frac{1}{n_{r}n_{f}}\right)^{k}\sum_{rf_{1}}\sum_{rf_{2}}...\sum_{rf_{k}}\|D_{rf_{k}}^{\frac{1}{2}}D_{rf_{1}}^{\frac{1}{2}}\|F\|D_{rf_{1}}^{\frac{1}{2}}D_{rf_{2}}^{\frac{1}{2}}\|F...\|D_{rf_{k-1}}^{\frac{1}{2}}D_{rf_{k}}^{\frac{1}{2}}\|F,$$

$$= d\left(\frac{1}{n_{r}n_{f}}\right)^{k}\sum_{rf_{1}}\sum_{rf_{2}}...\sum_{rf_{k}}S_{rf_{k},rf_{1}}S_{rf_{1},rf_{2}}...S_{rf_{k-1},rf_{k}},$$

$$= d\left(\frac{1}{n_{r}n_{f}}\right)^{k}\operatorname{Tr}(S^{k}),$$

$$\leq d^{2}\left(\frac{1}{n_{r}n_{f}}\right)^{k}\lambda_{1}(S)^{k}.$$

$$(33)$$

Therefore, we say that

$$\operatorname{Tr}[D^k] \le d^2 \left(\frac{1}{n_r n_f}\right)^k \lambda_1(S)^k,\tag{34}$$

and we can have that

$$\frac{(n_{r}n_{f})^{k} \operatorname{Tr}[D^{k}]}{d^{2}\sigma^{k}} \leq \frac{(n_{r}n_{f})^{k} \operatorname{Tr}[D^{k}]}{d^{2}\lambda_{1}(S)^{k}} \max_{i \in D_{r}j \in D_{f}} \lambda_{1} \left(\frac{C_{r}^{\frac{1}{2}}}{C_{r}^{\frac{1}{2}} + C_{f}^{\frac{1}{2}}} H_{i} + \frac{C_{f}^{\frac{1}{2}}}{C_{r}^{\frac{1}{2}} + C_{f}^{\frac{1}{2}}} H_{j}\right)^{k},$$

$$\leq \sum_{a_{r} \in D_{r}} \sum_{a_{f} \in D_{f}} \operatorname{Tr}\left[\left(\frac{C_{r}^{\frac{1}{2}}}{C_{r}^{\frac{1}{2}} + C_{f}^{\frac{1}{2}}} H_{a_{r}} + \frac{C_{f}^{\frac{1}{2}}}{C_{r}^{\frac{1}{2}} + C_{f}^{\frac{1}{2}}} H_{a_{f}}\right)^{k}\right].$$
(35)

Therefore, we can conclude that

$$\operatorname{Tr}[N_{k}] \geq \frac{1}{d} \frac{1}{n_{f} n_{r}} \frac{1}{2^{k}} \left(C_{r}^{\frac{1}{2}} + C_{f}^{\frac{1}{2}}\right)^{2k} \left(\frac{(n_{r} n_{f})^{k} \operatorname{Tr}[D^{k}]}{d^{2} \sigma^{k}}\right)^{2},$$

$$\geq \frac{1}{d^{5}} \frac{1}{n_{f} n_{r}} \frac{1}{2^{k}} \left(C_{r}^{\frac{1}{2}} + C_{f}^{\frac{1}{2}}\right)^{2k} \frac{(n_{r} n_{f})^{2k} \lambda_{1}(D)^{2k}}{\sigma^{2k}},$$

$$\geq \frac{1}{d^{5}} \frac{1}{n_{f} n_{r}} \frac{1}{2^{k}} \left(C_{r}^{\frac{1}{2}} + C_{f}^{\frac{1}{2}}\right)^{2k} \frac{(n_{r} n_{f})^{2k} \lambda_{1}(D)^{2k}}{\sigma^{2k}}.$$
(36)

Lastly, we can see that whether the trace diverge or not depend on those term with power of k. Therefore, by rearranging and plug in the definition of the coefficient into those terms, we can have that

$$\lambda_1(D) \ge \frac{\sqrt{2}\sigma}{\eta} \left( (1 - \alpha)n_f \left( \frac{n_r}{B} - 1 \right)^{\frac{1}{2}} + \alpha n_r \left( \frac{n_f}{B} - 1 \right)^{\frac{1}{2}} \right)^{-1},\tag{37}$$

which is the condition for diverging behavior

### 4.5 Proof of theorem 2.3

 **Theorem 4.9** ((Restate) Matching lower bound.). Suppose  $\lambda_{\max}(D)$  and  $\sigma$  satisfy

$$\lambda_{\max}(D) \le \frac{2\sigma}{\eta C_r' \left(\sigma + n_f \left(\frac{n_r}{B} - 1\right)\right)},\tag{38}$$

where  $C_r' = \sqrt{C_r}/(\sqrt{C_r} + \sqrt{C_f})$  (with  $C_r, C_f$  from Definition 2). Then there exists a choice of PSD Hessians  $\{H_i\}$  for the retain and forget sets such that the unlearning update converges (i.e.  $\lim_{k\to\infty} \mathbb{E}\|w_k\|^2 = 0$ ) under those Hessians.

*Proof.* We prove by construction in the following manner. We construct the retain set by setting  $H_i=m\cdot e_1e_1^T \ \ \forall \ i \in [\frac{\sigma}{n_f}]. \ (m=(C_r')^{-1}\frac{\lambda_1(D)n_r}{\frac{\sigma}{n_f}} \ \text{and} \ C_r'=\frac{C_r^{\frac{1}{2}}}{C_r^{\frac{1}{2}}+C_f^{\frac{1}{2}}} \ \text{and} \ \text{the definition of} \ C_r \ \text{and}$ 

 $C_f$  are mentioned in definition 2.) Otherwise, we set the Hessian to be zero matrix. For the forget set, we set all matrix to be zero matrix.

We first verify that the eigenvalue of mix-Hessian is indeed the assigned value  $\lambda_1(D)$ .

$$D = \frac{1}{n_r n_f} \sum_{r,f} C'_r H_r + C'_f H_f = \frac{1}{n_r n_f} \sum_{r,f} C'_r (C'_r)^{-1} \frac{\lambda_1(D) n_r}{\frac{\sigma}{n_f}} e_1 e_1^T = \lambda_1(D) e_1 e_1^T, \quad (39)$$

and we have that the construction indeed have the corresponding mix-Hessian eigenvalue.

We know verify that the coherence measure is of the assigned value  $\sigma$ . We first note that the element of the coherence matrix is:

$$S_{rf,r'f'} = \sqrt{\text{Tr}[(C'_r H_r + C'_f H_f)(C'_r H_{r'} + C'_f H_{f'})]} = C'_r m = \frac{\lambda_1(D)n_r}{\frac{\sigma}{n_f}}, \ \forall r, r' \in \left[\frac{\sigma}{n_f}\right].$$
(40)

else it is zero. We know that there is  $n_f \cdot \frac{\sigma}{n_f} = \sigma$  nonzero elements for each row and column. We note that we will also need to divide the coherence matrix by  $\max_{rf} D_{rf} = \max_{rf} C'_r H_r + C'_f H_f = \frac{\lambda_1(D)n_r}{\frac{\sigma}{n_f}}$ . Finally, each element is 1 after this division, and we can get the eigenvalue of the matrix to be  $\sigma$  and verify that the construction is valid.

Now, we note that in our construction, we have each step  $J_i$  to commute to each other since every matrix involved is diagonal, so we can focus on one step to calculate the condition that lead to diverging or converging and since we only intentionally set our matrix to be one dimensional, we can study behavior on only one axis  $e_1$  by plugging in the above as follows:

$$e_{1}E[J_{1}J_{1}]e_{1} = e_{1}[I - 2\eta(1-\alpha)H_{R} + \eta^{2}(1-\alpha)^{2}H_{r}^{2} + \eta^{2}(1-\alpha)^{2}\sum_{r}H_{r}^{2}]e_{1},$$

$$= 1 - 2\eta(1-\alpha)(C_{r}')^{-1}\lambda_{1}(D) + \eta^{2}(1-\alpha)^{2}(C_{r}')^{-2}\lambda_{1}(D)^{2} + \frac{(C_{r}')^{-2}}{\sigma}\lambda_{1}(D)^{2}\eta^{2}(1-\alpha)^{2}n_{f}(\frac{n_{r}}{B} - 1).$$
(41)

As we want to study the converging behavior, we want the above to be smaller than 1 to have repetitive multiplication lead to converging.

$$\begin{array}{ll}
1018 \\
1019 & 1 - 2\eta(1 - \alpha)(C_r')^{-1}\lambda_1(D) + \eta^2(1 - \alpha)^2(C_r')^{-2}\lambda_1(D)^2 + \frac{(C_r')^{-2}}{\sigma}\lambda_1(D)^2\eta^2(1 - \alpha)^2n_f(\frac{n_r}{B} - 1) \le 1, \\
1020 & \Longrightarrow 2 \ge \eta(1 - \alpha)(C_r')^{-1}\lambda_1(D)(1 + \frac{n_f}{\sigma}(\frac{n_r}{B} - 1)), \\
1022 & \Longrightarrow 2 \ge \frac{\eta}{\sigma}(1 - \alpha)(C_r')^{-1}\lambda_1(D)(\sigma + n_f(\frac{n_r}{B} - 1)), \\
1023 & \Longrightarrow \lambda_1(D) \le \frac{2\sigma}{\eta}C_r'\Big((1 - \alpha)(\sigma + n_f(\frac{n_r}{B} - 1))^{-1}\Big).
\end{array}$$

$$(42)$$

#### 4.6 Proof of theorem 4.6

**Theorem 4.10** (Restate). Under the data model of Definition 5 and the two-layer ReLU CNN defined above, suppose the network is trained to near-zero training loss. Then with probability at least  $1-8\delta$  (over the random draw of the dataset), the largest eigenvalue of the coherence matrix S for the retain/forget split satisfies

$$\lambda_{\max}(S) \leq \mathcal{O}\left(n_r \, n_f \, d\sigma^2 \left[ (\sqrt{C_r'} + \sqrt{C_f'})^2 \, (\text{SNR})^2 + (C_r' + C_f') \right] \right),$$
 (43)

$$\max_{rf} \lambda_{\max}(D_{rf}) \le \mathcal{O}((C_r' + C_f')(d\sigma^2(SNR)^2 + 1)), \tag{44}$$

where  $C'_r$  and  $C'_f$  are the normalized retain/forget weight fractions as defined in Theorem 2.3. Consider division of two quantities and we can find that for small SNR limit and large SNR limit:

$$\lim_{SNR \to 0} \frac{\lambda_{\max}(S)^{upper}}{\max_{rf} D_{rf}^{upper}} = \mathcal{O}(n_r n_f) , \lim_{SNR \to \infty} \frac{\lambda_{\max}(S)^{upper}}{\max_{rf} D_{rf}^{upper}} = \mathcal{O}(n_r n_f (1 + \frac{2\sqrt{C_r' C_f'}}{C_r' + C_f'})).$$
(45)

*Proof.* We first calculate the gradient of one sample respective to one of the  $w_{i,r}$ .

$$\frac{\partial \ell(y_i \cdot f(W, x_i))}{\partial w_{i,r}} = \ell_i' \cdot \frac{j}{m} \cdot (\mathbf{1}_{\{\langle w_{j,r}, y_i \cdot \boldsymbol{\mu} \rangle > 0\}} \cdot \boldsymbol{\mu} + \mathbf{1}_{\{\langle w_{j,r}, \xi_i \rangle > 0\}} \cdot y_i \cdot \xi_i). \tag{46}$$

There are several index in the above equation (i.e., j and r) which we use to take derivative with respect to a specific feature weight vector. We will continue to use this notation for future calculation. Now, we move to calculate the second derivative with respect to two different feature of weights for data i as follows:

$$\frac{\partial^{2}\ell(y_{i} \cdot f(W, x_{i}))}{\partial w_{j,r} \partial w_{j',r'}} = \ell_{i}^{"} \cdot \frac{jj'}{m^{2}} \cdot (\mathbf{1}_{\{\langle w_{j,r}, y_{i} \cdot \boldsymbol{\mu} \rangle > 0\}} \cdot \boldsymbol{\mu} + \mathbf{1}_{\{\langle w_{j,r}, \xi_{i} \rangle > 0\}} \cdot y_{i} \cdot \xi_{i}) (\mathbf{1}_{\{\langle w_{j',r'}, y_{i} \cdot \boldsymbol{\mu} \rangle > 0\}} \cdot \boldsymbol{\mu} + \mathbf{1}_{\{\langle w_{j',r'}, \xi_{i} \rangle > 0\}} \cdot y_{i} \cdot \xi_{i})^{T}.$$
(47)

The above is one block of the Hessian. In the following, we will simplify the notation for indicator function (derivative of ReLU) to  $\mathbf{1}_{j',r',y_i\cdot\mu}$  and  $\mathbf{1}_{j',r',\xi}$  to ease the heavy notation. To calculate the coherence matrix, we need to calculate trace of Hessian product for different sample,

$$\operatorname{Tr}[H_{i}H_{k}] = \sum_{j,j',r,r'} \operatorname{Tr}\left[\frac{\partial^{2}\ell(y_{i} \cdot f(W,x_{i}))}{\partial w_{j,r}\partial w_{j',r'}} \frac{\partial^{2}\ell(y_{k} \cdot f(W,x_{k}))}{\partial w_{j',r'}\partial w_{j,r}}\right], \\
= \frac{\ell_{i'}''\ell_{k}''}{m^{4}} \sum_{j,j',r,r'} (\mathbf{1}_{j,r,y_{k} \cdot \mu} \cdot \mu + \mathbf{1}_{j,r,\xi_{k}} \cdot y_{k} \cdot \xi_{k})^{T} (\mathbf{1}_{j,r,y_{i} \cdot \mu} \cdot \mu + \mathbf{1}_{j,r,\xi_{i}} \cdot y_{i} \cdot \xi_{i}) \\
(\mathbf{1}_{j',r',y_{i} \cdot \mu} \cdot \mu + \mathbf{1}_{j',r',\xi_{i}} \cdot y_{i} \cdot \xi_{i})^{T} (\mathbf{1}_{j',r',y_{k} \cdot \mu} \cdot \mu + \mathbf{1}_{j',r',\xi_{k}} \cdot y_{k} \cdot \xi_{k}), \\
= \frac{\ell_{i'}''\ell_{k}''}{m^{4}} \left( (\sum_{j,r} \mathbf{1}_{j,r,y_{k} \cdot \mu} \mathbf{1}_{j,r,y_{i} \cdot \mu}) \|\mu\|^{2}) + (\sum_{j,r} \mathbf{1}_{j,r,y_{k} \cdot \mu} \mathbf{1}_{j,r,\xi_{k}}) \mu^{T} \xi_{k} + (\sum_{j,r} \mathbf{1}_{j,r,\xi_{k}} \mathbf{1}_{j,r,\xi_{i}}) \xi_{k}^{T} \xi_{i}), \\
((\sum_{j,r} \mathbf{1}_{j',r',y_{k} \cdot \mu} \mathbf{1}_{j',r',y_{i} \cdot \mu}) \|\mu\|^{2}) + (\sum_{j',r'} \mathbf{1}_{j',r',y_{k} \cdot \mu} \mathbf{1}_{j',r',\xi_{k}}) \mu^{T} \xi_{k} + (\sum_{j',r'} \mathbf{1}_{j',r',\xi_{k}} \mathbf{1}_{j',r',\xi_{k}}) \xi_{k}^{T} \xi_{i}), \\
\leq 4 \frac{\ell_{i'}''\ell_{k}''}{m^{2}} (\|\mu\|^{2} + |\mu^{T} \xi_{k}| + |\mu^{T} \xi_{i}| + |\xi_{i}^{T} \xi_{k}|)^{2}, \\
\leq \frac{4}{m^{2}} (\|\mu\|^{2} + |\mu^{T} \xi_{k}| + |\mu^{T} \xi_{i}| + |\xi_{i}^{T} \xi_{k}|)^{2}.$$

We now analyze each term in the coherence matrix.

$$S_{r_{1}f_{1'},r_{2}f_{2'}} = \sqrt{\operatorname{Tr}((C'_{r}H_{r_{1}} + C'_{f}H_{f_{1'}})(C'_{r}H_{r_{2}} + C'_{f}H_{f_{2'}}))},$$

$$= \sqrt{\operatorname{Tr}[C'_{r}H_{r_{1}}H_{r_{2}}] + \operatorname{Tr}[C'_{r}C'_{f}H_{r_{1}}H_{f_{2'}}] + \operatorname{Tr}[C'_{r}C'_{f}H_{f_{1'}}H_{r_{2}}] + \operatorname{Tr}[C'_{f}H_{f_{1'}}H_{r_{2'}}]},$$

$$\leq \sqrt{\operatorname{Tr}[C'_{r}H_{r_{1}}H_{r_{2}}]} + \sqrt{\operatorname{Tr}[C'_{r}C'_{f}H_{r_{1}}H_{f_{2'}}]} + \sqrt{\operatorname{Tr}[C'_{r}C'_{f}H_{f_{1'}}H_{r_{2}}]} + \sqrt{\operatorname{Tr}[C'_{r}C'_{f}H_{f_{1'}}H_{r_{2'}}]},$$

$$(49)$$

where the  $C'_r$  and  $C'_f$  are respectively the normalized coefficient mentioned in the previous section.

As our goat is to estimate the largest eigenvalue of the coherence matrix and its relation between different variables in the design. To estimate the largest eigenvalue, we incur  $\epsilon$ -net that is used random matrix theory

$$\lambda_1 = \sup_{\|x\|=1} \langle x, Sx \rangle. \tag{50}$$

For one vector x, we can write the expression as summation:

$$\langle x, Sx \rangle = \sum_{r_{1}f_{1'}, r_{2}f_{2'}} S_{r_{1}f_{1'}, r_{2}f_{2'}} x_{r_{1}f_{1'}} x_{r_{2}f_{2'}},$$

$$\leq \sum_{r_{1}f_{1'}, r_{2}f_{2'}} (\sqrt{\operatorname{Tr}[C_{r}^{\prime 2}H_{r_{1}}H_{r_{2}}]} + \sqrt{\operatorname{Tr}[C_{r}^{\prime}C_{f}^{\prime}H_{r_{1}}H_{f_{2'}}]} + \sqrt{\operatorname{Tr}[C_{r}^{\prime}C_{f}^{\prime}H_{f_{1'}}H_{r_{2}}]} + \sqrt{\operatorname{Tr}[C_{r}^{\prime 2}H_{f_{1'}}H_{r_{2}}]} + \sqrt{\operatorname{Tr}[C_{r}^{\prime 2}H_{f_{1'}}H_{r_{2'}}]} + \sqrt{\operatorname{Tr}[C_{r}^{\prime 2}H_{f_{1'}}H_{r_{2'}}H_{$$

We can estimate the above through the random matrix theory and upper bound the largest eigenvalue through the elementwise calculation that we set up and use the tail bound for each random variable to provide relationship between controlled variable and the resulting largest eigenvalue. We first separate the discussion into several cases. First case, when we have four different samples  $r_1, r_2, f_1', f_2'$ , we can have that

1135
1136
$$(\sqrt{\text{Tr}[C_{r}^{\prime 2}H_{r_{1}}H_{r_{2}}]} + \sqrt{\text{Tr}[C_{r}^{\prime}C_{f}^{\prime}H_{r_{1}}H_{f_{2}^{\prime}}]} + \sqrt{\text{Tr}[C_{r}^{\prime}C_{f}^{\prime}H_{f_{1}^{\prime}}H_{r_{2}}]} + \sqrt{\text{Tr}[C_{f}^{\prime 2}H_{f_{1}^{\prime}}H_{r_{2}^{\prime}}]})x_{r_{1}f_{1}^{\prime}}x_{r_{2}f_{2}^{\prime}},$$
1137
1138
$$\leq (C_{r}^{\prime}\frac{2}{m}(\|\boldsymbol{\mu}\|^{2} + |\boldsymbol{\mu}^{T}\xi_{r1}| + |\boldsymbol{\mu}^{T}\xi_{r2}| + |\xi_{r_{1}}^{T}\xi_{r2}|) + \sqrt{C_{r}^{\prime}C_{f}^{\prime}}\frac{2}{m}(\|\boldsymbol{\mu}\|^{2} + |\boldsymbol{\mu}^{T}\xi_{r1}| + |\boldsymbol{\mu}^{T}\xi_{f2^{\prime}}| + |\xi_{r_{1}}^{T}\xi_{f2^{\prime}}|),$$
1139
$$+ \sqrt{C_{r}^{\prime}C_{f}^{\prime}}\frac{2}{m}(\|\boldsymbol{\mu}\|^{2} + |\boldsymbol{\mu}^{T}\xi_{r2}| + |\boldsymbol{\mu}^{T}\xi_{f1^{\prime}}| + |\xi_{f1^{\prime}}^{T}\xi_{r2}|) + C_{f}^{\prime}\frac{2}{m}(\|\boldsymbol{\mu}\|^{2} + |\boldsymbol{\mu}^{T}\xi_{f2^{\prime}}| + |\boldsymbol{\mu}^{T}\xi_{f1^{\prime}}| + |\xi_{f1^{\prime}}^{T}\xi_{f2^{\prime}}|))x_{r_{1}f_{1}^{\prime}}x_{r_{2}f_{2^{\prime}}},$$
1141
$$\leq (\sqrt{C_{r}^{\prime}}\|\boldsymbol{\mu}\| + \sqrt{C_{f}^{\prime}}\|\boldsymbol{\mu}\| + \sqrt{C_{f}^{\prime}}\|\xi_{r1}\| + \sqrt{C_{f}^{\prime}}\|\xi_{f1^{\prime}}\|)(\sqrt{C_{r}^{\prime}}\|\boldsymbol{\mu}\| + \sqrt{C_{f}^{\prime}}\|\boldsymbol{\mu}\| + \sqrt{C_{f}^{\prime}}\|\xi_{f2^{\prime}}\|)x_{r_{1}f_{1}^{\prime}}x_{r_{2}f_{2^{\prime}}}.$$
1143

Our aiming in the above is to establish relationship between different variables used in the CNN network. In the above, we can see that we can upper bound the eigenvalue by the cross product of the vector  $v_{rf} = \sqrt{C'_r} \|\boldsymbol{\mu}\| + \sqrt{C'_f} \|\boldsymbol{\mu}\| + \sqrt{C'_r} \|\boldsymbol{\xi}_{r1}\| + \sqrt{C'_f} \|\boldsymbol{\xi}_{f1'}\|$  since the coherence matrix is upper bound elementwise by the vector. (i.e.,  $\lambda_1(S) \leq \lambda_1(vv^T) = \|v^Tv\|^2$ ) and this turns the estimation of the eigenvalue into estimation of the magnitude of the vector.

Now, we analyze the  $v^T v$ ,

1152
1153
$$v^{T}v = \sum_{rf} (\sqrt{C'_{r}} \|\boldsymbol{\mu}\| + \sqrt{C'_{f}} \|\boldsymbol{\xi}_{r1}\| + \sqrt{C'_{f}} \|\boldsymbol{\xi}_{f1'}\|) (\sqrt{C'_{r}} \|\boldsymbol{\mu}\| + \sqrt{C'_{f}} \|\boldsymbol{\mu}\| + \sqrt{C'_{f}} \|\boldsymbol{\xi}_{f1'}\|),$$
1154
$$= \sum_{rf} (\sqrt{C'_{r}} \|\boldsymbol{\mu}\| + \sqrt{C'_{f}} \|\boldsymbol{\mu}\|)^{2} + 2(\sqrt{C'_{r}} \|\boldsymbol{\mu}\| + \sqrt{C'_{f}} \|\boldsymbol{\mu}\|) (\sqrt{C'_{r}} \|\boldsymbol{\xi}_{r1}\| + \sqrt{C'_{f}} \|\boldsymbol{\xi}_{f1'}\|) + (\sqrt{C'_{r}} \|\boldsymbol{\xi}_{f1'}\|)^{2},$$
1157
$$= n_{r} n_{f} (\sqrt{C'_{r}} \|\boldsymbol{\mu}\| + \sqrt{C'_{f}} \|\boldsymbol{\mu}\|)^{2} + 2(\sqrt{C'_{r}} \|\boldsymbol{\mu}\| + \sqrt{C'_{f}} \|\boldsymbol{\mu}\|) \sum_{rf} (\sqrt{C'_{r}} \|\boldsymbol{\xi}_{r1}\| + \sqrt{C'_{f}} \|\boldsymbol{\xi}_{f1'}\|) +$$
1159
$$= \sum_{rf} (C'_{r} \|\boldsymbol{\xi}_{r1}\|^{2} + C'_{f} \|\boldsymbol{\xi}_{f1'}\|^{2} + \sqrt{C'_{r}} C'_{f} \|\boldsymbol{\xi}_{r1}\| \|\boldsymbol{\xi}_{f1'}\|).$$
1160
$$\sum_{rf} (C'_{r} \|\boldsymbol{\xi}_{r1}\|^{2} + C'_{f} \|\boldsymbol{\xi}_{f1'}\|^{2} + \sqrt{C'_{r}} C'_{f} \|\boldsymbol{\xi}_{r1}\| \|\boldsymbol{\xi}_{f1'}\|).$$
1162
$$= \sum_{rf} (C'_{r} \|\boldsymbol{\xi}_{r1}\|^{2} + C'_{f} \|\boldsymbol{\xi}_{f1'}\|^{2} + \sqrt{C'_{r}} C'_{f} \|\boldsymbol{\xi}_{r1}\| \|\boldsymbol{\xi}_{f1'}\|).$$
1163

We analyze different terms as follows:

$$2(\sqrt{C'_r}\|\boldsymbol{\mu}\| + \sqrt{C'_f}\|\boldsymbol{\mu}\|) \sum_{rf} (\sqrt{C'_r}\|\xi_{r1}\| + \sqrt{C'_f}\|\xi_{f1'}\|) = 2(\sqrt{C'_r}\|\boldsymbol{\mu}\| + \sqrt{C'_f}\|\boldsymbol{\mu}\|)(n_f \sum_r \sqrt{C'_r}\|\xi_{r1}\| + n_r \sum_f \sqrt{C'_f}\|\xi_{f1'}\|).$$
(54)

We know that  $\|\xi_{r1}\|$ ,  $\|\xi_{f1'}\|$  are chi-distribution which is also sub-exponential distribution. We can utilize the tail bound for summation of the sub-exponential random variables to obtain high probability bound on the summation. We can have that with probability  $2\delta$ ,

$$2(\sqrt{C'_r}\|\boldsymbol{\mu}\| + \sqrt{C'_f}\|\boldsymbol{\mu}\|) \sum_{rf} (\sqrt{C'_r}\|\xi_{r1}\| + \sqrt{C'_f}\|\xi_{f1'}\|),$$

$$\leq 2(\sqrt{C'_r}\|\boldsymbol{\mu}\| + \sqrt{C'_f}\|\boldsymbol{\mu}\|) (n_r n_f \sqrt{C'_r} \sigma \sqrt{d} + n_f n_r \sqrt{C'_f} \sigma \sqrt{d} + n_f \sqrt{\frac{n_r \sigma^2}{C_1} \log(\frac{2}{\delta})} + n_r \sqrt{\frac{n_f \sigma^2}{C_1} \log(\frac{2}{\delta})}).$$
(55)

Now, we move to the next chi-square distribution terms  $C'_r \sum \|\xi_{r_1}\|^2$ ,  $C'_f \sum \|\xi_{f'_1}\|^2$ . By using the lemma 4.6, we can have that with probability  $1 - \delta$ ,

$$C'_r \sum_{rf} \|\xi_{r_1}\|^2 \le C'_r (n_f n_r \sigma^2 d + n_f \sqrt{\frac{n_r \sigma^4 d}{C_2} \log(\frac{2}{\delta})}).$$
 (56)

and so is the  $C'_f \sum \|\xi_{f'_1}\|^2$ ,

$$C_f' \sum_{rf} \|\xi_{r_1}\|^2 \le C_f' (n_f n_r \sigma^2 d + n_r \sqrt{\frac{n_f \sigma^4 d}{C_2} \log(\frac{2}{\delta})}).$$
 (57)

 The term  $\sum_{rf} \sqrt{C'_r C'_f} \|\xi_{r1}\| \|\xi_{f1'}\|$  can also be dealt with in the same manner,

$$\sqrt{C_r'C_f'} \sum_{rf} \|\xi_{r1}\| \|\xi_{f1'}\| \le \sqrt{C_r'C_f'} \left(\sum_r \|\xi_{r1}\|\right) \left(\sum_f \|\xi_{f1'}\|\right),$$

$$\le \sqrt{C_r'C_f'} \left(n_r \sqrt{2}\sigma \frac{\Gamma((d+1)/2)}{\Gamma(d/2)} + \sqrt{\frac{n_r \sigma^2}{C_1} \log(\frac{2}{\delta})}\right) \left(n_f \sqrt{2}\sigma \frac{\Gamma((d+1)/2)}{\Gamma(d/2)} + \sqrt{\frac{n_f \sigma^2}{C_1} \log(\frac{2}{\delta})}\right).$$
(58)

To simplify the analysis, we only keep terms with magnitude at least  $n_f n_r$ . We will reach that with probability  $1-6\delta$ 

$$\lambda_{1}(S) \leq \mathcal{O}\left(n_{f}n_{r}\left((\sqrt{C'_{r}} + \sqrt{C'_{f}})^{2}\|\boldsymbol{\mu}\|^{2} + 2\sqrt{2}(\sqrt{C'_{r}} + \sqrt{C'_{f}})^{2}\|\boldsymbol{\mu}\|\sigma\left(\frac{\Gamma((d+1)/2)}{\Gamma(d/2)}\right) + (C'_{r} + C'_{f})\sigma^{2}d + 2\sigma^{2}\sqrt{C'_{r}C'_{f}}\left(\frac{\Gamma((d+1)/2)}{\Gamma(d/2)}\right)^{2}\right)\right).$$
(59)

To see how signal noise ratio (SNR =  $\frac{\|\mu\|}{\sigma\sqrt{d}}$ ) interact with the right hand side, we extract a factor  $\sigma^2 d$  from all terms involved:

$$\begin{split} \lambda_{1}(S) &\leq \mathcal{O}\Big(n_{f}n_{r}\sigma^{2}d\Big((\sqrt{C'_{r}}+\sqrt{C'_{f}})^{2}(\text{SNR})^{2} + \frac{2\sqrt{2}}{\sqrt{d}}(\sqrt{C'_{r}}+\sqrt{C'_{f}})^{2}\big(\frac{\Gamma((d+1)/2)}{\Gamma(d/2)}\big)(\text{SNR}) \\ &+ (C'_{r}+C'_{f}) + \frac{2}{d}\sqrt{C'_{r}C'_{f}}\big(\frac{\Gamma((d+1)/2)}{\Gamma(d/2)}\big)^{2}\big)\Big), \\ &\leq \mathcal{O}\Big(n_{f}n_{r}\sigma^{2}d\big((\sqrt{C'_{r}}+\sqrt{C'_{f}})^{2}(\text{SNR})^{2} + (C'_{r}+C'_{f})\big)\Big). \end{split}$$
(60)

where in the last equation, we omit terms with d in the denominator as it tends to be large when we consider larger network.

For the second part of the proof, we know that  $H_i$  have block structures as follows:

$$\frac{\partial^{2}\ell(y_{i} \cdot f(W, x_{i}))}{\partial w_{j,r} \partial w_{j',r'}} = \ell_{i}^{"} \cdot \frac{jj'}{m^{2}} \cdot (\mathbf{1}_{\{\langle w_{j,r}, y_{i} \cdot \boldsymbol{\mu} \rangle > 0\}} \cdot \boldsymbol{\mu} + \mathbf{1}_{\{\langle w_{j,r}, \xi_{i} \rangle > 0\}} \cdot y_{i} \cdot \xi_{i}) (\mathbf{1}_{\{\langle w_{j',r'}, y_{i} \cdot \boldsymbol{\mu} \rangle > 0\}} \cdot \boldsymbol{\mu} + \mathbf{1}_{\{\langle w_{j',r'}, \xi_{i} \rangle > 0\}} \cdot y_{i} \cdot \xi_{i})^{T}.$$
(61)

We can see that the whole  $H_i$  matrix can be regarded as outer product of vector  $vv^T$  where we have  $v_{jr}$  being

$$v_{jr} = \frac{l_i''j}{m} (\mathbf{1}_{\{\langle w_{j,r}, y_i \cdot \mu \rangle > 0\}} \cdot \mu + \mathbf{1}_{\{\langle w_{j,r}, \xi_i \rangle > 0\}} \cdot y_i \cdot \xi_i).$$
 (62)

 We can immediately know that the eigenvalue of the  $H_i$  will be upper bounded by  $v^Tv$  as follows:

$$\lambda_{\max}(H_{i}) \leq v^{T} v = \frac{l_{i}^{"}}{m^{2}} \sum_{jr} (\mathbf{1}_{\{\langle w_{j,r}, y_{i} \cdot \boldsymbol{\mu} \rangle > 0\}} \cdot \boldsymbol{\mu} + \mathbf{1}_{\{\langle w_{j,r}, \xi_{i} \rangle > 0\}} \cdot y_{i} \cdot \xi_{i})^{2},$$

$$\leq 2 \frac{l_{i}^{"}}{m^{2}} \sum_{jr} \mathbf{1}_{\{\langle w_{j,r}, y_{i} \cdot \boldsymbol{\mu} \rangle > 0\}} \|\boldsymbol{\mu}\|^{2} + \mathbf{1}_{\{\langle w_{j,r}, \xi_{i} \rangle > 0\}} \|\xi_{i}\|^{2},$$

$$\leq 2 \frac{l_{i}^{"}}{m^{2}} \sum_{jr} \|\boldsymbol{\mu}\|^{2} + \|\xi_{i}\|^{2},$$

$$\leq 2 \frac{1}{m^{2}} \sum_{jr} \|\boldsymbol{\mu}\|^{2} + \|\xi_{i}\|^{2},$$

$$= C(\|\boldsymbol{\mu}\|^{2} + \|\xi_{i}\|^{2}),$$
(63)

where we use C to encompass all constants.

To bound the  $\max_{rf} \lambda_{\max}(D_{rf}) = \lambda_{\max}(C'_r H_r + C'_f H_f)$ , we can use the following:

$$\lambda_{\max}(D_{rf}) = \lambda_{\max}(C_r'H_r + C_f'H_f) \le C_r'\lambda_{\max}(H_r) + C_f'\lambda_{\max}(H_f). \tag{64}$$

Then for any  $\delta \in (0,1)$ , with probability at least  $1-\delta$ , we can upper bound the the  $H_r$  with the following ( $\|\xi_i\|$  is subexponential):

$$\max_{1 \le i \le n_r} C(\|\boldsymbol{\mu}\|^2 + \|\xi_i\|^2) \le C\left(\|\boldsymbol{\mu}\|^2 + \sigma^2 \left[d + 2\sqrt{d\log\frac{n_r}{\delta}} + 2\log\frac{n_r}{\delta}\right]\right), \quad (65)$$

$$\le \mathcal{O}(\|\boldsymbol{\mu}\|^2 + \sigma^2 d).$$

Similarly, we can have the bound on  $H_f$  which is of same order and jointly we can have that with probability  $1-8\delta$ 

$$\lambda_{\max}(D_{rf}) \le \mathcal{O}((C'_r + C'_f)(\|\boldsymbol{\mu}\|^2 + \sigma^2 d)) = \mathcal{O}((C'_r + C'_f)\sigma^2 d(SNR^2 + 1))$$
 (66)

Last is the division and take the limit and we can have the following:

$$\lim_{\text{SNR}\to 0} \frac{\lambda_{\text{max}}(S)^{\text{upper}}}{\max_{rf} D_{rf}^{\text{upper}}} = \mathcal{O}(n_r n_f) , \lim_{\text{SNR}\to \infty} \frac{\lambda_{\text{max}}(S)^{\text{upper}}}{\max_{rf} D_{rf}^{\text{upper}}} = \mathcal{O}(n_r n_f (1 + \frac{2\sqrt{C_r' C_f'}}{C_r' + C_f'})).$$
(67)