

528 **A Some concentration results for uniform random variables**

529 In this section, we state some concentration results that are useful for the theoretical analysis in
 530 Section 3. Let $\tilde{\mathbf{x}}, \mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathcal{B}_{\mathbf{0}, \sqrt{d+2}})$ be i.i.d. samples from the uniform distribution over
 531 the Euclidean norm ball of radius $\sqrt{d+2}$ in \mathbb{R}^d . Let

$$Z_n = \min_{\mathbf{x} \in \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}} \|\tilde{\mathbf{x}} - \mathbf{x}\|^2. \quad (11)$$

532 If $n = 1$, $\mathbb{E}Z_1$ is the sum of the variance of each coordinate of $\text{Unif}(\mathcal{B}_{\mathbf{0}, \sqrt{d+2}})$. Therefore, $\mathbb{E}Z_n$
 533 provides a generalized measure of concentration. Intuitively, $\mathbb{E}Z_n \rightarrow 0$ as $n \rightarrow \infty$. The proposition
 534 below provides an upper bound on the rate of convergence.

535 **Lemma A.1** (Nearest Neighbor concentration). *Given the assumptions above*

$$\mathbb{E}Z_n \lesssim d^2 \left[\frac{\log(n^{1/d})}{n} \right]^{1/d}, \quad (12)$$

536 where \lesssim means inequality up to an universal constant independent of d and n .

537 *Proof.* Define

$$\begin{aligned} \mathcal{E}_1 &= \{Z_n \leq \delta^2\}, \\ \mathcal{E}_2 &= \{\delta \leq \sqrt{d+2} - \|\tilde{\mathbf{x}}\|\}. \end{aligned} \quad (13)$$

538 We will compute two probabilities $\mathbb{P}(\mathcal{E}_1|\mathcal{E}_2)$ and $\mathbb{P}(\mathcal{E}_2)$ that will be useful latter.

$$\begin{aligned} \mathbb{P}(\mathcal{E}_1^c|\mathcal{E}_2) &= \mathbb{P}(Z_n \geq \delta^2|\mathcal{E}_2) = \mathbb{P}(\|\tilde{\mathbf{x}} - \mathbf{x}_i\| \geq \delta, \forall i|\mathcal{E}_2), \\ &= \mathbb{E}_{\tilde{\mathbf{x}}} \mathbb{P}(\|\tilde{\mathbf{x}} - \mathbf{x}_i\| \geq \delta|\mathcal{E}_2, \tilde{\mathbf{x}})^n = \mathbb{E}_{\tilde{\mathbf{x}}}(1 - \mathbb{P}(\|\tilde{\mathbf{x}} - \mathbf{x}_i\| \leq \delta|\mathcal{E}_2, \tilde{\mathbf{x}}))^n, \\ &= \mathbb{E}_{\tilde{\mathbf{x}}} \left[1 - \frac{\text{Vol}(\mathcal{B}_{\tilde{\mathbf{x}}, \delta})}{\text{Vol}(\mathcal{B}_{\mathbf{0}, \sqrt{d+2}})} \right]^n = \left[1 - \left(\frac{\delta}{\sqrt{d+2}} \right)^d \right]^n, \\ &\leq \exp \left[-n \left(\frac{\delta}{\sqrt{d+2}} \right)^d \right]. \end{aligned} \quad (14)$$

539 Next, we compute $\mathbb{P}(\mathcal{E}_2)$

$$\mathbb{P}(\mathcal{E}_2) = \mathbb{P}(\|\tilde{\mathbf{x}}\| \leq \sqrt{d+2} - \delta) = \left(\frac{\sqrt{d+2} - \delta}{\sqrt{d+2}} \right)^d = \left(1 - \frac{\delta}{\sqrt{d+2}} \right)^d. \quad (15)$$

540 We use \mathcal{E}_1 and \mathcal{E}_2 to compute the following upper bound

$$\begin{aligned} \mathbb{E}Z_n &= \mathbb{E}(Z_n|\mathcal{E}_1 \cap \mathcal{E}_2)\mathbb{P}(\mathcal{E}_1 \cap \mathcal{E}_2) + \mathbb{E}(Z_n|(\mathcal{E}_1 \cap \mathcal{E}_2)^c)\mathbb{P}((\mathcal{E}_1 \cap \mathcal{E}_2)^c), \\ &\leq \delta^2 + (2\sqrt{d+2})^2 (1 - \mathbb{P}(\mathcal{E}_1 \cap \mathcal{E}_2)), \\ &= \delta^2 + 4(d+2) [1 - \mathbb{P}(\mathcal{E}_1|\mathcal{E}_2)\mathbb{P}(\mathcal{E}_2)]. \end{aligned} \quad (16)$$

541 To find an upper bound for $\mathbb{E}Z_n$, we need to find an upper bound for $1 - \mathbb{P}(\mathcal{E}_1|\mathcal{E}_2)\mathbb{P}(\mathcal{E}_2)$.

$$\begin{aligned} 1 - \mathbb{P}(\mathcal{E}_1|\mathcal{E}_2)\mathbb{P}(\mathcal{E}_2) &= 1 - [1 - \mathbb{P}(\mathcal{E}_1^c|\mathcal{E}_2)]\mathbb{P}(\mathcal{E}_2), \\ &= 1 - \mathbb{P}(\mathcal{E}_2) + \mathbb{P}(\mathcal{E}_1^c|\mathcal{E}_2)\mathbb{P}(\mathcal{E}_2), \\ &\leq 1 - \mathbb{P}(\mathcal{E}_2) + \mathbb{P}(\mathcal{E}_1^c|\mathcal{E}_2). \end{aligned} \quad (17)$$

542 Now choose $\delta = \sqrt{d+2}n^{-1/d} [\log(n^{1/d})]^{1/d}$.

$$\mathbb{P}(\mathcal{E}_1^c|\mathcal{E}_2) \leq \exp \left[-n \left(\frac{\delta}{\sqrt{d+2}} \right)^d \right] = \exp \left[-n n^{-1} \log(n^{1/d}) \right] = n^{-1/d}, \quad (18)$$

543 and

$$\mathbb{P}(\mathcal{E}_2) = \left(1 - \frac{\delta}{\sqrt{d+2}}\right)^d \geq 1 - d \frac{\delta}{\sqrt{d+2}} = 1 - dn^{-1/d} \left[\log(n^{1/d})\right]^{1/d}. \quad (19)$$

544 Thus

$$1 - \mathbb{P}(\mathcal{E}_1|\mathcal{E}_2)\mathbb{P}(\mathcal{E}_2) \leq 1 - 1 + dn^{-1/d} \left[\log(n^{1/d})\right]^{1/d} + n^{-1/d} \lesssim dn^{-1/d} \left[\log(n^{1/d})\right]^{1/d}. \quad (20)$$

545 Combining everything together, we get

$$\begin{aligned} \mathbb{E}Z_n &\leq (d+2)n^{-2/d} \left[\log(n^{1/d})\right]^{2/d} + 4(d+2) \times dn^{-1/d} \left[\log(n^{1/d})\right]^{1/d}, \\ &\lesssim d^2 n^{-1/d} \left[\log(n^{1/d})\right]^{1/d}, \\ &= d^2 \left[\frac{\log(n^{1/d})}{n}\right]^{1/d}. \end{aligned} \quad (21)$$

546 This completes the proof. □

547 **Proposition A.2** ([47] Corollary 6.20). *Let $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathcal{B}_{\mathbf{0}, \sqrt{d+2}})$ for $i = 1, \dots, n$ be uniformly*
 548 *distributed over a ball of radius B in \mathbb{R}^d centered at $\mathbf{0}$. Let*

$$\boldsymbol{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$$

549 *be the sample covariance matrix. Then*

$$\mathbb{P}(\|\boldsymbol{\Sigma}_n - \mathbf{I}\|_{\text{op}} > \varepsilon) \leq 2d \exp\left[-\frac{n\varepsilon^2}{2(d+2)(1+\varepsilon)}\right].$$

550 **B Proof of Theorem 3.3**

551 In this section, we present the proof of Theorem 3.3. In Section B.1, we provide the detail of the
 552 decomposition of the risk into T_1 and T_2 . Then in Section B.2 we compute an upper bound for T_1 ,
 553 and compute an upper bound for T_2 in Section B.3. Finally, we combine everything together in
 554 Section B.4 and completes the proof.

555 **B.1 Decomposition of the test risk**

$$\begin{aligned}
 & \mathbb{E} [f^{\text{ResMem}}(\tilde{\mathbf{x}}) - f_*(\tilde{\mathbf{x}})]^2 = \mathbb{E} [f_n(\tilde{\mathbf{x}}) + r_n(\tilde{\mathbf{x}}) - f_*(\tilde{\mathbf{x}})]^2, \\
 & = \mathbb{E} [f_n(\tilde{\mathbf{x}}) - f_*(\tilde{\mathbf{x}}) - f_n(\tilde{\mathbf{x}}_{(1)}) + f_*(\tilde{\mathbf{x}}_{(1)})]^2, \\
 & = \mathbb{E} [f_n(\tilde{\mathbf{x}}) - f_\infty(\tilde{\mathbf{x}}) + f_\infty(\tilde{\mathbf{x}}) - f_*(\tilde{\mathbf{x}}) - f_n(\tilde{\mathbf{x}}_{(1)}) + f_\infty(\tilde{\mathbf{x}}_{(1)}) - f_\infty(\tilde{\mathbf{x}}_{(1)}) + f_*(\tilde{\mathbf{x}}_{(1)})]^2, \\
 & \leq 3 \times \underbrace{\mathbb{E}(f_n(\tilde{\mathbf{x}}) - f_\infty(\tilde{\mathbf{x}}))^2 + \mathbb{E}(f_n(\tilde{\mathbf{x}}_{(1)}) - f_\infty(\tilde{\mathbf{x}}_{(1)}))^2}_{T_1} + \underbrace{\mathbb{E}(f_\infty(\tilde{\mathbf{x}}) - f_*(\tilde{\mathbf{x}}) - f_\infty(\tilde{\mathbf{x}}_{(1)}) + f_*(\tilde{\mathbf{x}}_{(1)}))^2}_{T_2},
 \end{aligned} \tag{22}$$

556 where in the last inequality, we used the fact that $(a + b + c)^2 < 3(a^2 + b^2 + c^2)$ for any $a, b, c \in \mathbb{R}$.

557 **B.2 Upper bound on T_1 .**

558 Since $\mathbb{P}_{\mathbf{x}} = \text{Unif}(\mathcal{B}_{0,B})$, we apply the bound $\|\tilde{\mathbf{x}}\|, \|\tilde{\mathbf{x}}_{(1)}\| \leq B$ to obtain

$$\begin{aligned}
 T_1 & = \mathbb{E}[f_n(\tilde{\mathbf{x}}) - f_\infty(\tilde{\mathbf{x}})]^2 + \mathbb{E}[f_n(\tilde{\mathbf{x}}_{(1)}) - f_\infty(\tilde{\mathbf{x}}_{(1)})]^2, \\
 & = \mathbb{E}\langle \boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty, \tilde{\mathbf{x}} \rangle^2 + \mathbb{E}\langle \boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty, \tilde{\mathbf{x}}_{(1)} \rangle^2, \\
 & \leq \mathbb{E}\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|^2 \|\tilde{\mathbf{x}}\|^2 + \mathbb{E}\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|^2 \|\tilde{\mathbf{x}}_{(1)}\|^2, \\
 & \leq 2B^2 \mathbb{E}\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|^2.
 \end{aligned} \tag{23}$$

559 As n gets large, the empirical covariance matrix $\boldsymbol{\Sigma}_n = \mathbf{X}^\top \mathbf{X} / n$ is concentrated around its mean
 560 \mathbf{I} . Let $\boldsymbol{\Delta}_n = \mathbf{I} - \boldsymbol{\Sigma}_n$ denote this deviation. For some $\varepsilon \in (0, 1)$, define the following ‘‘good event’’
 561 over the randomness in $\boldsymbol{\Sigma}_n$

$$\mathcal{A} = \{\|\boldsymbol{\Delta}_n\|_{\text{op}} < \varepsilon\}, \tag{24}$$

562 where $\|\boldsymbol{\Delta}_n\|_{\text{op}}$ denotes the operator norm of the deviation matrix. The high level idea of the proof is
 563 to condition on the event \mathcal{A} and deduce an upper bound of $\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|$ in terms of ε . Then, we use
 564 the fact that \mathcal{A} happens with high probability.

565 Recall that $\boldsymbol{\theta}_\infty = L\boldsymbol{\theta}_*$, and

$$\boldsymbol{\theta}_n = \underset{\|\boldsymbol{\theta}\| \leq L}{\text{argmin}} \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2. \tag{25}$$

566 Since $\mathbf{y} = \mathbf{X}\boldsymbol{\theta}_*$ by definition, the Lagrangian of the convex program above is

$$\mathcal{L}(\boldsymbol{\theta}, \lambda) = \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{X}\boldsymbol{\theta}_*\|^2 + \lambda(\|\boldsymbol{\theta}\|^2 - L). \tag{26}$$

567 The KKT condition suggests that the primal-dual optimal pair $(\boldsymbol{\theta}_n, \lambda_n)$ is given by

$$\begin{aligned}
 \|\boldsymbol{\theta}_n\| & \leq L, \\
 \lambda_n & \geq 0, \\
 \lambda_n(\|\boldsymbol{\theta}_n\| - L) & = 0,
 \end{aligned} \tag{27}$$

568 and at optimality

$$\begin{aligned}
 \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_n, \lambda_n) = 0 & \iff \frac{2}{n} \mathbf{X}^\top \mathbf{X}(\boldsymbol{\theta}_n - \boldsymbol{\theta}_*) + 2\lambda_n \boldsymbol{\theta}_n = 0, \\
 & \iff \boldsymbol{\theta}_n = (\boldsymbol{\Sigma}_n + \lambda_n \mathbf{I})^{-1} \boldsymbol{\Sigma}_n \boldsymbol{\theta}_*.
 \end{aligned} \tag{28}$$

569 The complementary slackness condition $\lambda_n(\|\boldsymbol{\theta}_n\| - L) = 0$ suggests that either $\lambda_n = 0$ or $\|\boldsymbol{\theta}_n\| = L$.
570 But if $\lambda_n = 0$, the stationary condition $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \lambda) = 0$ would suggest that $\boldsymbol{\theta}_n = \boldsymbol{\Sigma}_n^{-1} \boldsymbol{\Sigma}_n \boldsymbol{\theta}_* =$
571 $\boldsymbol{\theta}_* \Rightarrow \|\boldsymbol{\theta}_n\| = 1 > L$, a contradiction. (Note that here $\boldsymbol{\Sigma}_n$ is invertible condition on the event \mathcal{A}).
572 Therefore, we must have $\|\boldsymbol{\theta}_n\| = L$. As a result, the primal and dual pair $(\boldsymbol{\theta}_n, \lambda_n)$ is determined by
573 the system of equations

$$\begin{cases} \boldsymbol{\theta}_n &= (\boldsymbol{\Sigma}_n + \lambda_n \mathbf{I})^{-1} \boldsymbol{\Sigma}_n \boldsymbol{\theta}_*, \\ \|\boldsymbol{\theta}_n\| &= L, \\ \lambda_n &> 0. \end{cases} \quad (29)$$

574 Next, we proceed to compute the deviation $\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|$.

$$\begin{aligned} \boldsymbol{\theta}_n &= [(\lambda_n + 1)\mathbf{I} - \boldsymbol{\Delta}_n]^{-1} \boldsymbol{\Sigma}_n \boldsymbol{\theta}_*, \\ &= (\lambda_n + 1)^{-1} \left[\mathbf{I} - \frac{\boldsymbol{\Delta}_n}{\lambda_n + 1} \right]^{-1} \boldsymbol{\Sigma}_n \boldsymbol{\theta}_*, \\ &= (\lambda_n + 1)^{-1} \left[\mathbf{I} + \sum_{k=1}^{\infty} \frac{\boldsymbol{\Delta}_n^k}{(\lambda_n + 1)^k} \right] (\mathbf{I} - \boldsymbol{\Delta}_n) \boldsymbol{\theta}_*, \\ &= (\lambda_n + 1)^{-1} \left[\mathbf{I} + \sum_{k=1}^{\infty} \frac{\boldsymbol{\Delta}_n^k}{(\lambda_n + 1)^k} - \boldsymbol{\Delta}_n - \sum_{k=1}^{\infty} \frac{\boldsymbol{\Delta}_n^{k+1}}{(\lambda_n + 1)^k} \right] \boldsymbol{\theta}_*, \\ &= (\lambda_n + 1)^{-1} \boldsymbol{\theta}_* + (\lambda_n + 1)^{-1} \boldsymbol{\Delta}_n \left[\sum_{k=1}^{\infty} \frac{\boldsymbol{\Delta}_n^{k-1}}{(\lambda_n + 1)^k} - \mathbf{I} - \sum_{k=1}^{\infty} \frac{\boldsymbol{\Delta}_n^k}{(\lambda_n + 1)^k} \right] \boldsymbol{\theta}_*, \\ &= (\lambda_n + 1)^{-1} \boldsymbol{\theta}_* + (\lambda_n + 1)^{-1} \boldsymbol{\Delta}_n \left[\sum_{k=1}^{\infty} \frac{\boldsymbol{\Delta}_n^{k-1} - \boldsymbol{\Delta}_n^k}{(\lambda_n + 1)^k} - \mathbf{I} \right] \boldsymbol{\theta}_*. \end{aligned} \quad (30)$$

575 Define

$$\mathbf{D}_n = \boldsymbol{\Delta}_n \left[\sum_{k=1}^{\infty} \frac{\boldsymbol{\Delta}_n^{k-1} - \boldsymbol{\Delta}_n^k}{(\lambda_n + 1)^k} - \mathbf{I} \right]. \quad (31)$$

576 Then $\boldsymbol{\theta}_n = (\lambda_n + 1)^{-1} \boldsymbol{\theta}_* + (\lambda_n + 1)^{-1} \mathbf{D}_n \boldsymbol{\theta}_*$, and

$$\begin{aligned} \|\mathbf{D}_n\| &\leq \|\boldsymbol{\Delta}_n\| \left[1 + \sum_{k=1}^{\infty} \frac{\|\boldsymbol{\Delta}_n\|^{k-1} + \|\boldsymbol{\Delta}_n\|^k}{(\lambda_n + 1)^k} \right], \\ &\leq \varepsilon \left[1 + 2(1 + \lambda_n)^{-1} \sum_{k=1}^{\infty} \left(\frac{\varepsilon}{1 + \lambda_n} \right)^k \right], \\ &= \varepsilon \left(1 + \frac{2}{1 + \lambda_n} \frac{1}{1 - \frac{\varepsilon}{1 + \lambda_n}} \right) \leq 3\varepsilon. \end{aligned} \quad (32)$$

577 Therefore

$$\begin{aligned} L &= \|\boldsymbol{\theta}_n\|^2 = (\lambda_n + 1)^{-2} + (\lambda_n + 1)^{-2} \boldsymbol{\theta}_*^\top \mathbf{D}_n^\top \mathbf{D}_n \boldsymbol{\theta}_* + 2(\lambda_n + 1)^{-2} \boldsymbol{\theta}_*^\top \mathbf{D}_n \boldsymbol{\theta}_*, \\ &\Rightarrow (\lambda_n + 1)^2 L^2 = 1 + \delta_n, \quad \delta_n = \boldsymbol{\theta}_*^\top \mathbf{D}_n^\top \mathbf{D}_n \boldsymbol{\theta}_* + 2\boldsymbol{\theta}_*^\top \mathbf{D}_n \boldsymbol{\theta}_*. \end{aligned} \quad (33)$$

578 We can obtain the following bound for δ_n :

$$|\delta_n| \leq \|\boldsymbol{\theta}_*\|^2 \|\mathbf{D}_n\|^2 + 2\|\boldsymbol{\theta}_*\| \|\mathbf{D}_n\| \leq 9\varepsilon^2 + 6\varepsilon \leq 15\varepsilon. \quad (34)$$

579 Since $1 - \delta_n/2 \leq \sqrt{1 + \delta_n} \leq 1 + \delta_n/2$, and $|\delta_n| \leq 15\varepsilon$, we obtain

$$|(\lambda_n + 1)L - 1| \leq \frac{15\varepsilon}{2} \Rightarrow |L - (\lambda_n + 1)^{-1}| \leq \frac{15\varepsilon}{2} (\lambda_n + 1)^{-1} \leq \frac{15\varepsilon}{2}, \quad (35)$$

580 where the last inequality follows as we have $\lambda_n > 0$. Finally,

$$\begin{aligned}
\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty &= (\lambda_n + 1)^{-1} \boldsymbol{\theta}_* - L \boldsymbol{\theta}_* + (\lambda_n + 1)^{-1} \mathbf{D}_n \boldsymbol{\theta}_*, \\
\Rightarrow \|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|^2 &= [(1 + \lambda_n)^{-1} - L]^2 + (1 + \lambda_n)^{-2} \boldsymbol{\theta}_*^T \mathbf{D}_n^T \mathbf{D}_n \boldsymbol{\theta}_* + 2(\lambda_n + 1)^{-1} [(1 + \lambda_n)^{-1} - L] \boldsymbol{\theta}_*^T \mathbf{D}_n \boldsymbol{\theta}_*, \\
&\leq 64\varepsilon^2 + 9\varepsilon^2 + 45\varepsilon^2 = 118\varepsilon^2, \\
\Rightarrow \|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|^2 &\lesssim \varepsilon^2.
\end{aligned} \tag{36}$$

581 Combine the above result with Proposition A.2, we get that

$$\begin{aligned}
\mathbb{E}\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|^2 &= \mathbb{E}(\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|^2 | \mathcal{A}) \mathbb{P}(\mathcal{A}) + \mathbb{E}(\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|^2 | \mathcal{A}^c) \mathbb{P}(\mathcal{A}^c), \\
&\leq \varepsilon^2 + 4L^2 \times 4d \exp\left[-\frac{n\varepsilon^2}{2(d+2)(1+\varepsilon)}\right],
\end{aligned} \tag{37}$$

582 If we choose $\varepsilon = n^{-1/3}$, we get

$$\mathbb{E}\|\boldsymbol{\theta}_n - \boldsymbol{\theta}_\infty\|^2 \lesssim dL^2 n^{-2/3}, \tag{38}$$

583 which implies that

$$T_1 \lesssim d^2 L^2 n^{-2/3}. \tag{39}$$

584 **B.3 Upper bound on T_2 .**

585 Plugging in the formula for $f_\perp(\tilde{\mathbf{x}}) = f_*(\tilde{\mathbf{x}}) - f_\infty(\tilde{\mathbf{x}}) = \langle \tilde{\mathbf{x}}, \boldsymbol{\theta}_\perp \rangle$, we get

$$\begin{aligned}
T_2 &= \mathbb{E}[f_\perp(\tilde{\mathbf{x}}_{(1)}) - f_\perp(\tilde{\mathbf{x}})]^2, \\
&= \mathbb{E}\langle \boldsymbol{\theta}_\perp, \tilde{\mathbf{x}}_{(1)} - \tilde{\mathbf{x}} \rangle^2, \\
&\leq (1 - L)^2 \|\boldsymbol{\theta}_*\|^2 \mathbb{E}\|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}_{(1)}\|^2, \\
&= (1 - L)^2 \mathbb{E}\|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}_{(1)}\|^2,
\end{aligned} \tag{40}$$

586 where in the last inequality, we used the relation that $\boldsymbol{\theta}_\perp = (1 - L)\boldsymbol{\theta}_*$. Proposition A.1 suggests that

$$\mathbb{E}\|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}_{(1)}\|^2 \lesssim d^2 \left[\frac{\log(n^{1/d})}{n} \right]^{1/d}, \tag{41}$$

587 which implies

$$T_2 \lesssim d^2 (1 - L)^2 \left[\frac{\log(n^{1/d})}{n} \right]^{1/d}. \tag{42}$$

588 *Remark B.1* (Comparison with pure nearest neighbor and ERM). If we rely solely on nearest neighbor
589 method, the prediction error is

$$\mathbb{E}[f_*(\tilde{\mathbf{x}}) - f_*(\tilde{\mathbf{x}}_{(1)})]^2 = \mathbb{E}\langle \tilde{\mathbf{x}} - \tilde{\mathbf{x}}_{(1)}, \boldsymbol{\theta}_* \rangle^2 \leq \mathbb{E}\|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}_{(1)}\|^2. \tag{43}$$

590 On the other hand, if we solely rely on ERM, even with infinite sample, we get

$$\mathbb{E}[f_*(\tilde{\mathbf{x}}) - f_\infty(\tilde{\mathbf{x}})]^2 = \mathbb{E}\langle \tilde{\mathbf{x}}, \boldsymbol{\theta}_* - \boldsymbol{\theta}_\infty \rangle^2 \leq (1 - L)^2 \mathbb{E}\|\tilde{\mathbf{x}}\|^2. \tag{44}$$

591 We can see from the upper bound that ResMem takes advantage of both

- 592 • Projecting f_* onto f_∞ , so that the dependence on the prediction function is reduced from 1
593 to $(1 - L)^2$.
- 594 • Memorizing the residuals using nearest neighbor, so that the variance is reduced from $\mathbb{E}\|\tilde{\mathbf{x}}\|^2$
595 to $\mathbb{E}\|\tilde{\mathbf{x}}_{(1)} - \tilde{\mathbf{x}}\|^2$.

596 **B.4 Test loss for ResMem.**

597 If we combine the previous two parts together, we get

$$\mathbb{E} \left[\hat{f}(\tilde{\mathbf{x}}) - f_*(\tilde{\mathbf{x}}) \right]^2 \lesssim d^2 L^2 n^{-2/3} + d^2 (1-L)^2 \left[\frac{\log(n^{1/d})}{n} \right]^{1/d}. \quad (45)$$

598 This completes the proof of Theorem 3.3.

599 **C Additional CIFAR100 Results**

600 This section includes additional experiment results on applying ResMem to CIFAR100 dataset.

601 **C.1 Additional robustness results**

602 In addition to the results already presented in Section 4.2, we also evaluate ResMem performance for
 603 each architecture in CIFAR-ResNet{8, 14, 20, 32, 44, 56} and each subset (10%, 20%, ..., 100%) of
 604 CIFAR100 training data. We use the same training hyperparameter and the ResMem hyperparameter
 605 as described in Section 4.2. Generally, we see that ResMem yields larger improvement over the
 baseline DeepNet when the network is small and dataset is large.

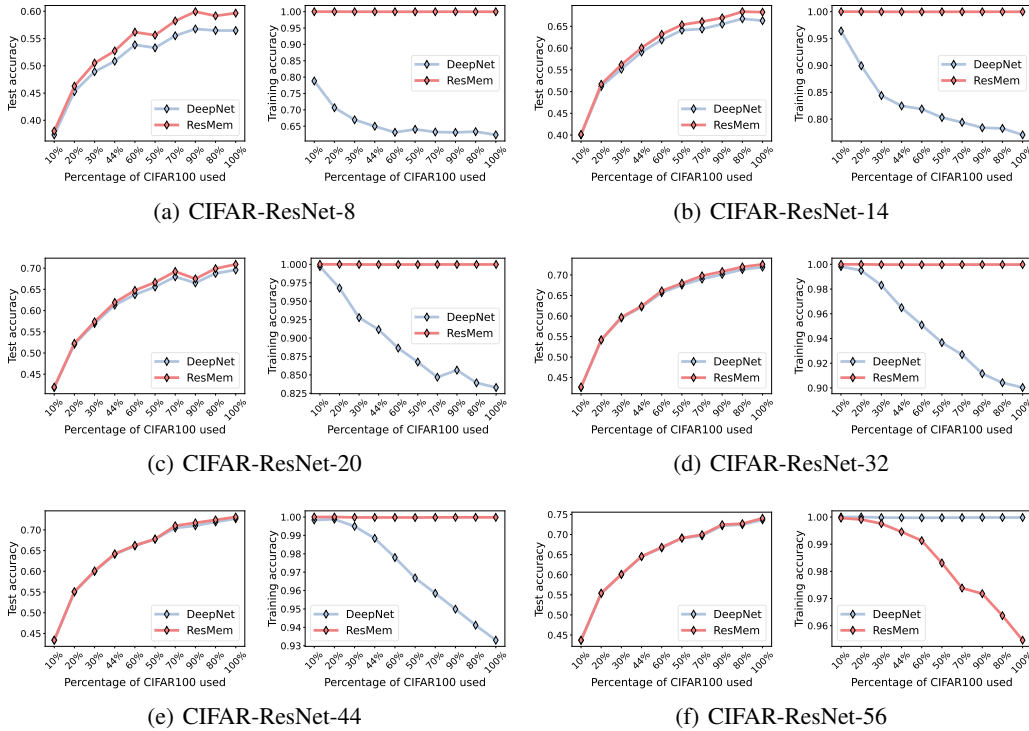


Figure 4: Test(left)/Training (right) accuracy for different sample sizes.

606

607 **C.2 Sensitivity analysis for CIFAR100**

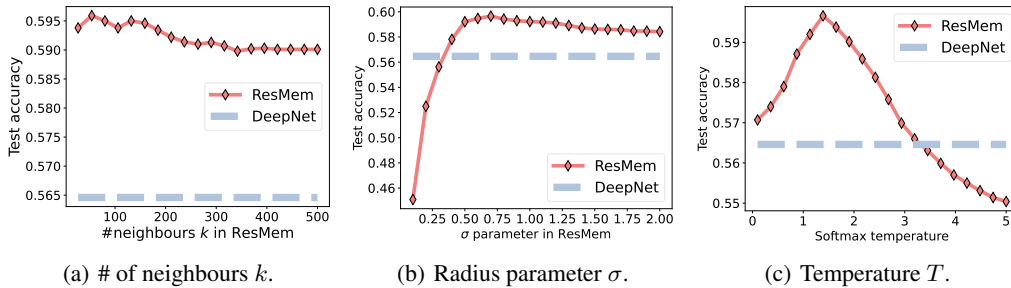


Figure 5: Sensitivity analysis of ResMem hyperparameters. The y -axis represents the CIFAR100 test accuracies, and the x -axis represents the sweeping of respective hyperparameters.

608 **Varying locality parameter k and σ .** We vary the number of neighbours from $k = 27$ to $k = 500$.
 609 We find that ResMem test accuracy is relatively stable across the choice of the number of neighbours
 610 (cf. Figure 5(a)). The trend of the curve suggests that as $k \rightarrow \infty$, the ResMem test accuracy seems to
 611 be converging to a constant level. For σ , we explored different values of $\sigma \in (0.1, 2.0)$. We observe
 612 that the test accuracy has a unimodal shape as a function of σ , suggesting that there is an optimal
 613 choice of σ (cf. Figure 5(b)).

614 **Varying temperature T and connection to distillation.** We tried $T = 0.1$ to $T = 5$, and also
 615 identified a unimodal shape for the test accuracy (Figure 5(c)). The fact that we can use different
 616 temperatures for (a) training the network and (b) constructing the k -NN predictor reminds us of the
 617 well-established knowledge distillation procedure [28]. In knowledge distillation, we first use one
 618 model (the teacher network) to generate targets at a higher temperature, and then train a second model
 619 (the student network) using the *combination* of the true labels and the output of the first network.

620 ResMem operates in a reversed direction: Here we have a second model (kNN) that learns the
 621 *difference* between true labels and the output of the first model. In both cases, we can tune the
 622 temperature of the first model to control how much information is retained. This connection offers an
 623 alternative perspective that regards ResMem as a “dual procedure” to knowledge distillation.

624 D ResMem on ImageNet

625 This section includes additional experiment results on applying ResMem to ImageNet dataset.

626 **ImageNet.** In addition to CIFAR100, we also evaluate the performance of ResMem on ImageNet [42].
 627 We employ a family of pre-trained MobileNet-V2 models [44] from Keras², with
 628 varying widths controlled by a multiplier a . For ResMem, we again use the second last layer of
 629 DeepNet as a 1280-dimensional embedding of an image and rely on the ℓ_2 distance between the
 630 embeddings for nearest neighbor search (Step 3, Section 4.1). We specify the ResMem parameter of
 631 (k, σ, T) in the table below. We repeat the experiment over several MobileNet-V2 architectures, with
 632 MobileNet-V2-a0.35 being the smallest model and MobileNet-V2-a1.3 being the largest one.

Table 1: Test accuracy for ResMem and baseline deep network for ImageNet data.

Architecture	ResMem param.			Test accuracy	
	k	σ	T	DeepNet	ResMem
MobileNet-V2-a0.35	10	0.6	0.4	60.2%	61.2%
MobileNet-V2-a0.5	10	0.6	0.4	65.3%	66.1%
MobileNet-V2-a0.75	10	0.8	0.6	69.6%	70.1%
MobileNet-V2-a1.0	20	0.4	0.4	71.3%	71.8%
MobileNet-V2-a1.3	30	0.4	0.4	74.7%	75.1%

633 We can see that (c.f. Table 1) ResMem boosts the test accuracy by 1% on the smallest model and by
 634 0.4% on the largest model.

635 E Additional details of NLP experiments

636 The Decoder-Only model used in our experiments is essentially the normal Encoder-Decoder archi-
 637 tecture with Encoder and Cross-Attention removed. We pretrained both the T5-small and T5-base
 638 model on C4 [41] dataset with auto-regressive language modeling task for 1,000,000 steps, with
 639 dropout rate of 0.1 and batch size of 128. The learning rate for the first 10,000 steps is fixed to 0.01
 640 and the rest steps follow a square root decay schedule.

641 During the inference for retrieval key, query embeddings and residuals, we ensured every token has
 642 at least 64 preceding context by adopting a sliding window strategy, where a window of 256 token
 643 slides from the beginning to the end on each of the articles, with a stride of $256 - 64 = 192$.

²<https://keras.io/api/applications/mobilenet/>

644 For residuals, we only stored the top 128 residuals measured by the absolute magnitude, as the
645 residual vector is as large as T5 vocabulary size (i.e., 32128), and storing all 32128 residuals for each
646 token is too demanding for storage. However, when weight-combining the residuals, we zero filled
647 the missing residuals so that all the residual vectors have 32128 elements.