

A Additional Discussions

A.1 Limitations

In this work, we focus on the optimal membership inference adversary. We study this because of how it serves as an upper bound for all other attacks and because of how it yields interpretable and fundamental theoretical results. The optimal membership inference adversary has full knowledge of the learning model’s output distributions when the data point of interest is a member or non-member of the training dataset. In practice, the adversary rarely has such full knowledge, and the learning model’s output distributions have to be approximated using shadow models [13], or the entire attack has to be simplified, such as with a loss threshold [9, 14]. Our study does not analyze how our results are affected by the non-optimality of these more practical attacks.

A.2 Ethical Considerations

It is the hope of the authors that by more clearly exposing the link between membership inference vulnerability and generalization performance, researchers can make informed decisions about how to achieve the best trade-off they can for their application. That said, by studying the performance of optimal membership inference attacks, it is possible that this work will call attention to vulnerabilities in existing model architectures which may be exploited. Furthermore, in settings where privacy is absolutely crucial, such as in medical applications, additional care should be taken to guard privacy beyond the guarantees of this work.

B Proofs

B.1 Proof of Proposition 3.1

We first present the proof of the form of the optimal membership inference adversary given in Proposition 3.1

Proof of Proposition 3.1 Conditioned on $m = 0$, we have that (\mathbf{x}_0, y_0) is drawn from \mathcal{D} . Conditioned on $m = 1$, we have that (\mathbf{x}_0, y_0) is an element chosen randomly from S , whose elements are themselves drawn from \mathcal{D} . Thus, in both the $m = 0$ and $m = 1$ cases, \mathbf{x}_0 has the same distribution. We thus have:

$$\begin{aligned} A^* &= \arg \max_A \text{Adv}(A) \\ &= \arg \max_A \mathbb{P}(A((\mathbf{x}_0, \hat{y}_0)) = 1 \mid m = 1) - \mathbb{P}(A((\mathbf{x}_0, \hat{y}_0)) = 1 \mid m = 0) \\ &= \arg \max_A \mathbb{E}_{\mathbf{x}_0} [\mathbb{P}(A((\mathbf{x}_0, \hat{y}_0)) = 1 \mid m = 1, \mathbf{x}_0) - \mathbb{P}(A((\mathbf{x}_0, \hat{y}_0)) = 1 \mid m = 0, \mathbf{x}_0)] \\ &= \arg \max_A \mathbb{E}_{\mathbf{x}_0} \left[\int_{\mathbb{R}} \mathbb{1}_{A(\mathbf{x}_0, \hat{y}_0)=1} (P(\hat{y}_0 \mid m = 1, \mathbf{x}_0) - P(\hat{y}_0 \mid m = 0, \mathbf{x}_0)) dP \right], \end{aligned}$$

where in the third line, the randomness over \mathbf{x}_0 is removed from the probability. To maximize the integral, we set $A(\mathbf{x}_0, \hat{y}_0) = 1$ if $P(\hat{y}_0 \mid m = 1, \mathbf{x}_0) - P(\hat{y}_0 \mid m = 0, \mathbf{x}_0) > 0$ and 0 otherwise. \square

B.2 Tools for Asymptotic Analysis

The following lemmas are used in the proofs of Theorems 3.2 and D.1. We begin with the following lemma, which is a generalized version of the Marchenko-Pastur theorem [34–36].

Lemma B.1. *Let $\mathbf{X}_n \in \mathbb{R}^{n \times p}$ be a sequence of random matrices with i.i.d. $\mathcal{N}(0, 1)$ entries. Consider the sample covariance matrix $\hat{\Sigma} = (1/n)\mathbf{X}_n^\top \mathbf{X}_n$. Let $\mathbf{C}_n \in \mathbb{R}^{p \times p}$ be a sequence of matrices such that $\text{Tr}(\mathbf{C}_n)$ is uniformly bounded with probability one. As $n, p \rightarrow \infty$ with $p/n = \gamma \in (0, \infty)$, it holds that almost surely,*

$$\text{Tr} \left(\mathbf{C}_n \left((\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} - g(-\lambda) \mathbf{I}_p \right) \right) \rightarrow 0, \quad \text{Tr} \left(\mathbf{C}_n \left((\hat{\Sigma} + \lambda \mathbf{I}_p)^{-2} - g'(-\lambda) \mathbf{I}_p \right) \right) \rightarrow 0$$

where $g(\lambda)$ is the Stieltjes transform of the Marchenko-Pastur law with parameter γ .

We use the following Lemma in computing the asymptotic distribution of the output.

Lemma B.2. *Let $\mathbf{y}_n \in \mathbb{R}^n$ be a sequence of i.i.d. $\mathcal{N}(0, \mathbf{I}_n)$ random vectors. Also, let $\mathbf{x}_n \in \mathbb{R}^n$ be a sequence of random vectors with spherically symmetric distribution such that $\|\mathbf{x}_n\|_2 \xrightarrow{a.s.} \sigma$. Further, assume that $\mathbf{x}_n, \mathbf{y}_n$ are independent. Then $\mathbf{x}_n^\top \mathbf{y}_n$ converges weakly to a zero mean gaussian with variance σ^2 .*

Proof. We can write

$$\mathbf{x}_n^\top \mathbf{y}_n = \|\mathbf{x}_n\|_2 \left(\frac{\mathbf{x}_n}{\|\mathbf{x}_n\|} \right)^\top \mathbf{y}_n = \|\mathbf{x}_n\|_2 \mathbf{u}_n^\top \mathbf{y}_n$$

where $\mathbf{u}_n \in S^{n-1}$ is uniformly distributed over the unit sphere and is independent from \mathbf{y}_n . Therefore, we can fix \mathbf{u}_n to be the first standard unit vector and the distribution of $\mathbf{x}_n^\top \mathbf{y}_n$ is the same as $\|\mathbf{x}_n\|_2 y_{n,1}$ where $y_{n,1} \sim \mathcal{N}(0, 1)$. Hence, using $\|\mathbf{x}_n\|_2 \xrightarrow{a.s.} \sigma$, we deduce the result. \square

B.3 Proof of Theorem 3.2

Proof of Theorem 3.2 Let $\mathbf{X}_{\bar{p}}$ denote the matrix formed by removing the first p columns from \mathbf{X} , and let $\beta_{\bar{p}}$ denote the vector formed by removing the first p elements from β . Recall that

$$\begin{aligned} (\hat{y}_0 \mid m = 0) &= \mathbf{x}_0^\top \mathbf{X}_p^\dagger (\mathbf{X}\beta + \epsilon) \\ &= \mathbf{x}_0^\top \mathbf{X}_p^\dagger (\mathbf{X}_p \beta_p + \eta) \end{aligned}$$

where $\eta = \mathbf{X}_{\bar{p}} \beta_{\bar{p}} + \epsilon \sim \mathcal{N}(0, (1 + \sigma^2 - \frac{p}{D}) \mathbf{I}_n)$. First note that the distributions of $\mathbf{X}_p^\top \mathbf{X}_p^\dagger \mathbf{x}_0$ are spherically symmetric and letting $\hat{\Sigma} \triangleq (1/n) \mathbf{X}_p^\top \mathbf{X}_p$ and \mathbf{P} to be orthogonal projection onto row space of \mathbf{X}_p we have

$$\begin{aligned} \frac{1}{D} \|\mathbf{X}_p^\top \mathbf{X}_p^\dagger \mathbf{x}_0\|_2^2 &= \frac{1}{D} \|\mathbf{P} \mathbf{x}_0\|_2^2 = \frac{1}{D} \lim_{\lambda \rightarrow 0} \mathbf{x}_0^\top (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \hat{\Sigma} \mathbf{x}_0 \\ &= \frac{1}{D} \|\mathbf{x}_0\|_2^2 - \frac{1}{D} \lim_{\lambda \rightarrow 0} \lambda \mathbf{x}_0^\top (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_0, \end{aligned}$$

and,

$$\begin{aligned} \|\mathbf{X}_p^\dagger \mathbf{x}_0\|_2^2 &= \frac{1}{n} \lim_{\lambda \rightarrow 0} \mathbf{x}_0^\top (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \hat{\Sigma} (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_0 \\ &= \frac{1}{n} \lim_{\lambda \rightarrow 0} \mathbf{x}_0^\top \left[(\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} - \lambda (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-2} \right] \mathbf{x}_0. \end{aligned}$$

Thus, using Lemma B.2, both $\frac{1}{D} \|\mathbf{X}_p^\top \mathbf{X}_p^\dagger \mathbf{x}_0\|_2^2$ and $\|\mathbf{X}_p^\dagger \mathbf{x}_0\|_2^2$ converge to a fixed limit as $n \rightarrow \infty$, almost surely. Therefore, using Lemma B.2, \hat{y}_0 converges weakly to a gaussian. Now, we compute its variance. Since η and β are both zero-mean independent Gaussians and are thus orthogonal in expectation, we have by the Pythagorean theorem:

$$\mathbb{E} \left[\hat{y}_0^2 \mid m = 0 \right] = \mathbb{E} \left[(\mathbf{x}_0^\top \mathbf{X}_p^\dagger \mathbf{X}_p \beta_p)^2 \right] + \mathbb{E} \left[(\mathbf{x}_0^\top \mathbf{X}_p^\dagger \eta)^2 \right].$$

We start with the first term.

Note that since, $p > n$, \mathbf{X}_p does not have linearly independent columns. Let $\mathbf{P} = \mathbf{X}_p^\dagger \mathbf{X}_p$. We have:

$$\begin{aligned}
\mathbb{E} \left[(\mathbf{x}_0^\top \mathbf{X}_p^\dagger \mathbf{X}_p \beta_p)^2 \right] &= \mathbb{E} \left[\text{Tr} \left(\mathbf{x}_0^\top \mathbf{P} \beta_p \beta_p^\top \mathbf{P}^\top \mathbf{x}_0 \right) \right] \\
&= \mathbb{E} \left[\text{Tr} \left(\beta_p \beta_p^\top \mathbf{P}^\top \mathbf{x}_0 \mathbf{x}_0^\top \mathbf{P} \right) \right] \\
&= \text{Tr} \left(\mathbb{E} \left[\beta_p \beta_p^\top \mathbf{P}^\top \mathbf{x}_0 \mathbf{x}_0^\top \mathbf{P} \right] \right) \\
&= \text{Tr} \left(\mathbb{E} \left[\beta_p \beta_p^\top \right] \mathbb{E} \left[\mathbf{P}^\top \mathbf{x}_0 \mathbf{x}_0^\top \mathbf{P} \right] \right) \\
&= \frac{1}{D} \text{Tr} \left(\mathbb{E} \left[\mathbf{P}^\top \mathbf{x}_0 \mathbf{x}_0^\top \mathbf{P} \right] \right) \\
&= \frac{1}{D} \mathbb{E} \left[\|\mathbf{P}^\top \mathbf{x}_0\|^2 \right] \\
&= \frac{1}{D} \frac{n}{p} \|\mathbf{x}_0\|^2.
\end{aligned}$$

In the last line, we use the same argument as in Section 2.2 of [23], using the facts that \mathbf{P} is the orthogonal projection to the row space of \mathbf{X}_p and that the Gaussian distribution is invariant to rotations.

We now consider the second term:

$$\begin{aligned}
\mathbb{E} \left[(\mathbf{x}_0^\top \mathbf{X}_p^\dagger \eta)^2 \right] &= \left(1 + \sigma^2 - \frac{p}{D} \right) \mathbf{x}_0^\top \mathbb{E} \left[\mathbf{X}_p^\dagger \mathbf{X}_p^{\dagger \top} \right] \mathbf{x}_0 \\
&= \left(1 + \sigma^2 - \frac{p}{D} \right) \mathbf{x}_0^\top \mathbb{E} \left[\left(\mathbf{X}_p^\top \mathbf{X}_p \right)^\dagger \mathbf{X}_p^\top \mathbf{X}_p \left(\mathbf{X}_p^\top \mathbf{X}_p \right)^{\dagger \top} \right] \mathbf{x}_0 \\
&= \left(1 + \sigma^2 - \frac{p}{D} \right) \mathbf{x}_0^\top \mathbb{E} \left[\left(\mathbf{X}_p^\top \mathbf{X}_p \right)^\dagger \right] \mathbf{x}_0.
\end{aligned}$$

where $\left(\mathbf{X}_p^\top \mathbf{X}_p \right)^\dagger$ has the generalized inverse Wishart distribution with expectation equal to $\mathbb{E} \left[\left(\mathbf{X}_p^\top \mathbf{X}_p \right)^\dagger \right] = \frac{n}{p} \frac{1}{p-n-1} \mathbf{I}_p$ (Theorem 2.1 in [37]). Thus, we have:

$$\mathbb{E} \left[(\mathbf{x}_0^\top \mathbf{X}_p^\dagger \eta)^2 \right] = \left(\frac{n}{p} \right) \left(\frac{1 + \sigma^2 - \frac{p}{D}}{p - n - 1} \right) \|\mathbf{x}_0\|^2$$

Adding this with the result for the first term gives the desired result. When $m = 1$, since we are in the overparameterized regime, \mathbf{X}_p is a fat matrix. Thus, the regressor memorizes the training data and the training error is equal to zero. \mathbf{x}_0 is part of training set, and so $\hat{y}_0 = \mathbf{x}_0^\top \beta + \epsilon$. Since $\beta \sim \mathcal{N} \left(0, \frac{1}{D} \mathbf{I}_p \right)$, we have that $\mathbf{x}_0^\top \beta \sim \mathcal{N} \left(0, \frac{1}{D} \|\mathbf{x}_0\|^2 \right)$. Since $\epsilon \sim \mathcal{N} \left(0, \sigma^2 \right)$, we have that $\hat{y}_0 = \mathbf{x}_0^\top \beta + \epsilon \sim \mathcal{N} \left(0, \frac{1}{D} \|\mathbf{x}_0\|^2 + \sigma^2 \right)$.

The probability distribution functions of the two Gaussians are then equal at $\pm\alpha$:

$$\begin{aligned} \frac{1}{\sigma_0\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{\alpha}{\sigma_0}\right)^2\right) &= \frac{1}{\sigma_1\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{\alpha}{\sigma_1}\right)^2\right) \\ \frac{\sigma_1}{\sigma_0} &= \exp\left(-\frac{1}{2}\left(\left(\frac{\alpha}{\sigma_1}\right)^2 - \left(\frac{\alpha}{\sigma_0}\right)^2\right)\right) \\ \frac{\sigma_1}{\sigma_0} &= \exp\left(-\frac{\alpha^2}{2}\frac{\sigma_0^2 - \sigma_1^2}{\sigma_0^2\sigma_1^2}\right) \\ \log\left(\frac{\sigma_1}{\sigma_0}\right) &= -\frac{\alpha^2}{2}\frac{\sigma_0^2 - \sigma_1^2}{\sigma_0^2\sigma_1^2} \\ \alpha &= \sqrt{\frac{2\sigma_0^2\sigma_1^2 \log\left(\frac{\sigma_1}{\sigma_0}\right)}{\sigma_1^2 - \sigma_0^2}} \\ \alpha &= \sqrt{\frac{\sigma_0^2\sigma_1^2 \log\left(\frac{\sigma_1}{\sigma_0}\right)}{\sigma_1^2 - \sigma_0^2}}. \end{aligned}$$

The membership advantage is then derived by writing out the probabilities in Definition 2.1 in terms of the Gaussian cumulative distribution functions, noting that the decision region switches at $\pm\alpha$. \square

Proof of Lemma 4.2. The lemma follows identically to Theorem 3.2 with an additional additive $\bar{\sigma}^2$ to σ_0^2 due to the noise added in the $m = 0$ case. The remainder follows by plugging in $p = D$ and applying Prop. 4.1 for the generalization error. \square

C Posterior Distributions in Non-Asymptotic Regime

Let $f_{a|b}$ denote the probability density function of a random variable a conditioned on b . The following lemma derives the non-asymptotic probability densities of the prediction output of minimum norm least squares, conditioned on the $m = 0$ and $m = 1$ events and the choice of test point \mathbf{x}_0 . For a matrix $\mathbf{X} \in \mathbb{R}^{n \times D}$ and $p \leq D$, let \mathbf{X}_p denote the submatrix of the first p columns of \mathbf{X} . For a vector $\mathbf{x} \in \mathbb{R}^D$, let $\mathbf{x}_p \in \mathbb{R}^p$ be defined accordingly.

Lemma C.1. *Let $\hat{\beta}$ denote the minimum norm least squares interpolator computed from a random design matrix $\mathbf{X} \in \mathbb{R}^{n \times D}$ and data \mathbf{y} . Conditioned on $n < p \leq D$ and on \mathbf{x}_0 , we have that $\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m = 1\} \sim \mathcal{N}(0, \sigma_1^2)$, where σ_1 is defined as in Theorem 3.2. Furthermore,*

$$\begin{aligned} & f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m=0\}}(\mathbf{x}) \\ &= D^{\frac{D}{2}} \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} \frac{\exp\left[-\frac{1}{2}\left[\left(\frac{\mathbf{x} - \mathbf{x}_{0,p}^\top \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-1} \mathbf{X} \beta\right)^2}{\sigma \|\mathbf{x}_{0,p}\| \|\mathbf{X}_p (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p\right]} + D\beta^\top \beta + \text{Tr}(\mathbf{X}^\top \mathbf{X})\right]}{\sigma(2\pi)^{\frac{nD+D+1}{2}} \|\mathbf{x}_{0,p}\| \|\mathbf{X}_p (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p\}} d\beta d\mathbf{X}. \end{aligned}$$

Remark C.2. While the density of $\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m = 0\}$ cannot be written in a closed form, one may easily sample according to it, by first, sampling random \mathbf{X} , β and then computing the minimum norm least squares interpolator.

Proof of Lemma C.1. Recall that conditioned on the design matrix \mathbf{X} and true coefficients β , the labels \mathbf{y} follow $\mathbf{y} \mid \mathbf{X}, \beta \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$. Then, the minimum norm least squares solution $\hat{\beta}$ using the first p features follows

$$\hat{\beta} \mid \mathbf{X}, \beta \sim \mathcal{N}\left(\mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-1} \mathbf{X} \beta, \sigma^2 \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p\right).$$

Hence, for the $m = 0$ case where a fresh point \mathbf{x}_0 is sampled, we have that the distribution of the model output conditioned on the design matrix \mathbf{X} and true coefficients β is

$$\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m = 0\} \sim \mathcal{N}\left(\mathbf{x}_{0,p}^\top \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-1} \mathbf{X} \beta, \sigma^2 \|\mathbf{x}_{0,p}\|_{\mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p}^2\right)$$

for $\|\mathbf{x}\|_{\mathbf{A}} := \sqrt{\mathbf{x}^\top \mathbf{A} \mathbf{x}}$ for any semidefinite matrix \mathbf{A} where we have additionally conditioned over any randomness in the choice of \mathbf{x}_0 .

In the $m = 1$ case, where \mathbf{x}_0 is sampled uniformly from the rows of \mathbf{X} , we have that $\mathbf{x}_{0,p}^\top \hat{\beta} = y_0 = \mathbf{x}_0^\top \beta + \epsilon$, the associated label for \mathbf{x}_0 since the linear regressor interpolates the training data. Hence

$$\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m = 1\} \equiv \mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m = 1\} \sim \mathcal{N}\left(0, \frac{\|\mathbf{x}_0\|^2}{D} + \sigma^2\right)$$

Let $f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m=0\}}$ denote the pdf of the random variable $\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m = 0\}$ and $f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m=1\}}$ be defined similarly. Let $f_{\mathbf{X}}$ denote the density of \mathbf{X} , a standard matrix-normal random variable, and let f_β denote the density of $\beta \sim \mathcal{N}(0, \frac{1}{D} \mathbf{I}_D)$. Then, we have that

$$\begin{aligned} & f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m=0\}}(x) \\ &= \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m=0\}} f_\beta f_{\mathbf{X}} d\beta d\mathbf{X} \\ &= D^{\frac{D}{2}} \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} \frac{\exp\left[-\frac{1}{2} \left[\left(\frac{x - \mathbf{x}_{0,p}^\top \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-1} \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| \mathbf{x}_p (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{x}_p} \right)^2 + D\beta^\top \beta + \text{Tr}(\mathbf{X}^\top \mathbf{X}) \right]\right]}{\sigma(2\pi)^{\frac{nD+D+1}{2}} \|\mathbf{x}_{0,p}\| \mathbf{x}_p (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{x}_p} d\beta d\mathbf{X}. \end{aligned}$$

□

Lemma C.3. Let $\hat{\beta}_\lambda = (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{y}$ denote ridge regularized least squares estimator computed from random design matrix $\mathbf{X} \in \mathbb{R}^{n \times D}$, data \mathbf{y} , and subset of first p features. Conditioned on the choice of test point \mathbf{x}_0 , we have that in the $m = 0$ case, where a fresh test point is drawn from the data distribution,

$$\begin{aligned} & f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m=0\}}(x) = \frac{D^{\frac{D}{2}}}{\sigma(2\pi)^{\frac{nD+D+1}{2}}} \\ & \times \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} \frac{\exp\left[-\frac{1}{2} \left[\left(\frac{x - \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{x}_p^\top \mathbf{x}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} \right)^2 + D\beta^\top \beta + \text{Tr}(\mathbf{X}^\top \mathbf{X}) \right]\right]}{\|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{x}_p^\top \mathbf{x}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} d\beta d\mathbf{X}. \end{aligned}$$

Furthermore, conditioned on $m = 1$ when \mathbf{x}_0 is a row of \mathbf{X} we have that

$$\begin{aligned} & f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m=1\}}(x) = \frac{D^{\frac{D}{2}}}{\sigma(2\pi)^{\frac{nD+1}{2}}} \\ & \times \int_{\mathbb{R}^{(n-1) \times D}} \int_{\mathbb{R}^D} \frac{\exp\left[-\frac{1}{2} \left[\left(\frac{x - \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{x}_p^\top \mathbf{x}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} \right)^2 + D\beta^\top \beta + \text{Tr}(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \right]\right]}{\|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{x}_p^\top \mathbf{x}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} d\beta d\tilde{\mathbf{X}}. \end{aligned}$$

where without loss of generality, we take \mathbf{x}_0 to be the first row of \mathbf{X} and $\tilde{\mathbf{X}}$ to denote the matrix of the final $n - 1$ rows of \mathbf{X} .

Remark C.4. As in the case of Lemma C.1, one can efficiently sample from the above distribution by first drawing a Gaussian random matrix \mathbf{X} , the Gaussian random vector β , the Bernoulli random variable m , and then either a new test point \mathbf{x}_0 or a row of \mathbf{X} and learning the ridge-regularized estimator $\hat{\beta}_\lambda$.

Proof of Lemma C.3. Note that conditioned on the design matrix \mathbf{X} and the true coefficients β , the labels \mathbf{y} follow $\mathbf{y} \mid \mathbf{X}, \beta \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$. Next, for

$$\hat{\beta}_\lambda := (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{y}$$

the λ -ridge regularized estimator, we have that

$$\hat{\beta}_\lambda \mid \mathbf{X}, \beta \sim \mathcal{N} \left((\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{X} \beta, \sigma^2 (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \right).$$

Hence,

$$\begin{aligned} \mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{X}, \mathbf{x}_0, \beta \\ \sim \mathcal{N} \left(\mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{X} \beta, \sigma^2 \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{x}_{0,p} \right), \end{aligned}$$

where we have additionally conditioned over any randomness in the choice of \mathbf{x}_0 .

In the $m = 0$ case, where \mathbf{x}_0 is a freshly drawn point independent of the data \mathbf{X} , we may marginalize to remove the conditioning. Let $f_{\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m=0\}}$ denote the probability density function of the random variable $\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m = 0\}$ and $f_{\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{X}, \mathbf{x}_0, \{m=1\}}$ be defined similarly. Let $f_{\mathbf{X}}$ denote the density of \mathbf{X} , a standard matrix-normal random variable, and let f_β denote the density of $\beta \sim \mathcal{N}(0, \frac{1}{D} \mathbf{I}_D)$. Then we have that

$$f_{\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{x}_0, \{m=0\}}(x) = \int_{\mathbb{R}^n \times D} \int_{\mathbb{R}^D} f_{\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{X}, \mathbf{x}_0, \{m=0\}} f_\beta f_{\mathbf{X}} d\beta d\mathbf{X}.$$

Thus,

$$\begin{aligned} f_{\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{x}_0, \{m=0\}}(x) &= \frac{D^{\frac{D}{2}}}{\sigma(2\pi)^{\frac{nD+D+1}{2}}} \\ &\times \int_{\mathbb{R}^n \times D} \int_{\mathbb{R}^D} \frac{\exp \left[-\frac{1}{2} \left[\left(\frac{x - \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} \right)^2 + D\beta^\top \beta + \text{Tr}(\mathbf{X}^\top \mathbf{X}) \right] \right]}{\|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} d\beta d\mathbf{X}. \end{aligned}$$

In the $m = 1$ case, because \mathbf{x}_0 is a row of \mathbf{X} , we condition on \mathbf{x}_0 but not on the remaining rows of \mathbf{X} . Without loss of generality, let \mathbf{x}_0 be the first row of \mathbf{X} which can be done since \mathbf{x}_0 is selected uniformly and the rows of \mathbf{X} are independent and identically distributed. Let $\tilde{\mathbf{X}} \in \mathbb{R}^{(n-1) \times D}$ denote \mathbf{X} with its first row omitted such that $\mathbf{X} = [\mathbf{x}_0; \tilde{\mathbf{X}}]$. Following the same approach as the preceding marginalization, we have that

$$f_{\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{x}_0, \{m=1\}}(x) = \int_{\mathbb{R}^{(n-1) \times D}} \int_{\mathbb{R}^D} f_{\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{X}, \mathbf{x}_0, \{m=1\}} f_\beta f_{\tilde{\mathbf{X}}} d\beta d\tilde{\mathbf{X}}$$

Thus,

$$\begin{aligned} f_{\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda \mid \mathbf{x}_0, \{m=1\}}(x) &= \frac{D^{\frac{D}{2}}}{\sigma(2\pi)^{\frac{nD+1}{2}}} \\ &\times \int_{\mathbb{R}^{(n-1) \times D}} \int_{\mathbb{R}^D} \frac{\exp \left[-\frac{1}{2} \left[\left(\frac{x - \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} \right)^2 + D\beta^\top \beta + \text{Tr}(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \right] \right]}{\|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} d\beta d\tilde{\mathbf{X}}. \end{aligned}$$

□

D Theoretical Results for Regularized Linear Regression

Theorem D.1. Membership advantage for Ridge-regularized linear regression. Consider the same setup as in Theorem 3.2 but now let $\hat{\beta}_\lambda = (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X}$ for some $\lambda > 0$. Then, as $n, p, D \rightarrow \infty$ such that $\frac{p}{n} \rightarrow \gamma \in (1, \infty)$, we have:

$$\begin{aligned} (\hat{y}_0 \mid m = 0, \mathbf{x}_0) &\sim \mathcal{N}(0, \sigma_{0,\lambda}^2), \\ (\hat{y}_0 \mid m = 1, \mathbf{x}_0) &\sim \mathcal{N}(0, \sigma_{1,\lambda}^2), \end{aligned}$$

where

$$\begin{aligned}\sigma_{0,\lambda}^2 &:= \frac{g'(-\lambda)\gamma}{(1+g(-\lambda)\gamma)^2} \left(\sigma^2 + 1 - \frac{p}{D} \right) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} + (1 - 2\lambda g(-\lambda) + \lambda^2 g'(-\lambda)) \frac{\|\mathbf{x}_{0,p}\|_2^2}{D} \\ \sigma_{1,\lambda}^2 &:= \left[\left(\frac{\lambda^2}{(\lambda + \gamma g(-\lambda))(\lambda + \gamma \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} g(-\lambda))} \right)^2 \gamma g'(-\lambda) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} \right] \left(\sigma^2 + 1 - \frac{p}{D} \right) \\ &\quad + \left(\frac{\gamma g(-\lambda) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p}}{1 + \gamma g(-\lambda) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p}} \right)^2 \left(\sigma^2 + \frac{\|\mathbf{x}_{0,p}\|_2^2}{D} \right) \\ &\quad + \left(1 - \frac{2\lambda g(-\lambda)}{1 + \frac{\gamma \|\mathbf{x}_{0,p}\|_2^2}{p} g(-\lambda)} + \frac{\lambda^2 g'(-\lambda)}{\left(1 + \frac{\gamma \|\mathbf{x}_{0,p}\|_2^2}{p} g(-\lambda)\right)^2} \right) \frac{\|\mathbf{x}_{0,p}\|_2^2}{D}, \\ g(-\lambda) &:= \frac{-(1-\gamma+\lambda) + \sqrt{(1-\gamma+\lambda)^2 + 4\gamma\lambda}}{2\gamma\lambda}.\end{aligned}$$

Furthermore, in the case when $\sigma_{1,\lambda} > \sigma_{0,\lambda}$ and defining:

$$\alpha_\lambda = \sqrt{\frac{\sigma_{0,\lambda}^2 \sigma_{1,\lambda}^2 \log\left(\frac{\sigma_{1,\lambda}^2}{\sigma_{0,\lambda}^2}\right)}{\sigma_{1,\lambda}^2 - \sigma_{0,\lambda}^2}},$$

the optimal membership inference advantage is then:

$$\text{Adv}(A_\lambda^*) = \mathbb{E}_{\mathbf{x}_0} \left[2 \left\{ \Phi\left(\frac{\alpha_\lambda}{\sigma_{0,\lambda}}\right) - \Phi\left(\frac{\alpha_\lambda}{\sigma_{1,\lambda}}\right) \right\} \right].$$

Remark D.2. The above result holds using the asymptotic distributions as $n, p, D \rightarrow \infty$. In Lemma C.3, we derive the non-asymptotic distributions for the predictions of the ridge-regularized least squares estimator, though they cannot be written in closed form.

Proof of Theorem D.1 Let the input be $\mathbf{x}_{0,p} \in \mathbb{R}^p$. Similar to the proof of theorem 3.2, we can write

$$\mathbf{X}\beta + \epsilon = \mathbf{X}_p\beta_p + \mathbf{X}_{\bar{p}}\beta_{\bar{p}} + \epsilon = \mathbf{X}_p\beta_p + \eta$$

where $\eta = \mathbf{X}_{\bar{p}}\beta_{\bar{p}} + \epsilon \sim \mathcal{N}(0, (1 + \sigma^2 - \frac{p}{D}) \mathbf{I}_n)$. Hence, we have

$$\hat{y}_0 = \mathbf{x}_{0,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top (\mathbf{X}_p\beta_p + \eta) = \mathbf{x}_{0,p}^\top \mathbf{X}_p^\top \left(\mathbf{X}_p \mathbf{X}_p^\top + n\lambda \mathbf{I}_p \right)^{-1} (\mathbf{X}_p\beta_p + \eta). \quad (7)$$

First note that in the case $m = 0$, we have

$$\begin{aligned}\hat{y}_0 &= \mathbf{x}_{0,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top (\mathbf{X}_p\beta_p + \eta) \\ &= \mathbf{x}_{0,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{X}_p\beta_p + \mathbf{x}_{0,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \eta.\end{aligned} \quad (8)$$

Letting the sample covariance matrix $\hat{\Sigma} \triangleq (1/n) \mathbf{X}_p^\top \mathbf{X}_p$, the first term in (8) can be written as

$$\begin{aligned}\mathbf{x}_{0,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{X}_p\beta_p &= \mathbf{x}_{0,p}^\top \left(\hat{\Sigma} + \lambda \mathbf{I}_p \right)^{-1} \hat{\Sigma} \beta_p \\ &= \mathbf{x}_{0,p}^\top \left(\hat{\Sigma} + \lambda \mathbf{I}_p \right)^{-1} \left(\hat{\Sigma} + \lambda \mathbf{I}_p - \lambda \mathbf{I}_p \right) \beta_p \\ &= \left(\mathbf{x}_{0,p} - \lambda \left(\hat{\Sigma} + \lambda \mathbf{I}_p \right)^{-1} \mathbf{x}_{0,p} \right)^\top \beta_p.\end{aligned}$$

Since $\beta_p \sim (0, \frac{1}{D}\mathbf{I}_D)$ using Lemma [B.2](#) this converges to a gaussian with zero mean and variance

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{D} \left\| \mathbf{x}_{0,p} - \lambda(\widehat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_{0,p} \right\|_2^2 &= \lim_{n \rightarrow \infty} \frac{1}{D} \left\{ \|\mathbf{x}_{0,p}\|_2^2 - 2\lambda[\mathbf{x}_{0,p}^\top (\widehat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_{0,p}] \right. \\ &\quad \left. + \lambda^2[\mathbf{x}_{0,p}^\top (\widehat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p)^{-2} \mathbf{x}_{0,p}] \right\} \\ &= \frac{\|\mathbf{x}_{0,p}\|_2^2}{D} (1 - 2\lambda g(-\lambda) + \lambda^2 g'(-\lambda)) \end{aligned}$$

where for the second equality, we have used the fact that using Lemma [B.1](#) by setting $\mathbf{C}_n = (1/n)\mathbf{x}_{0,p}\mathbf{x}_{0,p}^\top$, as $n \rightarrow \infty$, almost surely,

$$\frac{1}{n} \mathbf{x}_{0,p}^\top (\widehat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_{0,p} \rightarrow \frac{1}{n} \|\mathbf{x}_{0,p}\|_2^2 g(-\lambda), \quad \frac{1}{n} \mathbf{x}_{0,p}^\top (\widehat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p)^{-2} \mathbf{x}_{0,p} \rightarrow \frac{1}{n} \|\mathbf{x}_{0,p}\|_2^2 g'(-\lambda).$$

For the second term in [\(8\)](#), using the rotationally invariance of gaussian distribution, without loss of generality, we can let η to be $\mathbf{e}_1 \|\eta\|_2$, where \mathbf{e}_1 is the first standard unit vector. Now, note that we have

$$\|\eta\|_2 \mathbf{x}_{0,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{e}_1 = \|\eta\|_2 \mathbf{x}_{0,p}^\top \left(\mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top + \lambda n \mathbf{I}_p + \sum_{i=2}^n \mathbf{x}_{i,p} \mathbf{x}_{i,p}^\top \right)^{-1} \mathbf{x}_{1,p}$$

where $\mathbf{x}_{i,p}^\top \in \mathbb{R}^p$ is the i 'th row of \mathbf{X}_p . Letting $\mathbf{A}_\lambda \triangleq \lambda \mathbf{I}_p + \frac{1}{n} \sum_{i=2}^n \mathbf{x}_{i,p} \mathbf{x}_{i,p}^\top$, by using the Sherman-Morrison formula, we have

$$\begin{aligned} \|\eta\|_2 \mathbf{x}_{0,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{e}_1 &= \|\eta\|_2 \mathbf{x}_{0,p}^\top \left(\mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top + n\mathbf{A}_\lambda \right)^{-1} \mathbf{x}_{1,p} \\ &= \frac{\|\eta\|_2}{n} \mathbf{x}_{0,p}^\top \left(\mathbf{A}_\lambda^{-1} - \frac{\mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \right) \mathbf{x}_{1,p} \\ &= \frac{\|\eta\|_2}{n} \mathbf{x}_{0,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \left(1 - \frac{\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \right) \\ &= \|\eta\|_2 \frac{\mathbf{x}_{0,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}. \end{aligned}$$

Note that using Lemma [B.1](#) by setting $\mathbf{C}_n = (1/n)\mathbf{x}_{1,p}\mathbf{x}_{1,p}^\top$ and $\mathbf{C}_n = (1/n)\mathbf{x}_{0,p}\mathbf{x}_{0,p}^\top$, respectively, for $n, p \rightarrow \infty$, almost surely,

$$\frac{1}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \rightarrow \gamma g(-\lambda), \quad \frac{1}{n} \mathbf{x}_{0,p}^\top \mathbf{A}_\lambda^{-2} \mathbf{x}_{0,p} \rightarrow \frac{\|\mathbf{x}_{0,p}\|_2^2}{n} g'(-\lambda).$$

Thus, since $\mathbf{x}_{1,p} \sim (0, \mathbf{I}_p)$, using Lemma [B.2](#), $\|\eta\|_2 \mathbf{x}_{0,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{e}_1$ converges to a gaussian with mean zero and variance

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{\|\eta\|_2^2}{n^2 (1 + \gamma g(-\lambda))^2} \|\mathbf{A}_\lambda^{-1} \mathbf{x}_{0,p}\|_2^2 &= \frac{\|\eta\|_2^2 \|\mathbf{x}_{0,p}\|_2^2}{n^2 (1 + \gamma g(-\lambda))^2} g'(-\lambda) \\ &= \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} \frac{g'(-\lambda) \gamma}{1 + \gamma g(-\lambda))^2} \left(\sigma^2 + 1 - \frac{p}{D} \right). \end{aligned}$$

Hence, by independence of β_p and η , for $m = 0$, as $n \rightarrow \infty$, such that $p/n = \gamma$, the output \hat{y}_0 as in [\(7\)](#), converges in distribution to a gaussian with mean zero and variance

$$\frac{g'(-\lambda) \gamma}{(1 + g(-\lambda) \gamma)^2} \left(\sigma^2 + 1 - \frac{p}{D} \right) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} + (1 - 2\lambda g(-\lambda) + \lambda^2 g'(-\lambda)) \frac{\|\mathbf{x}_{0,p}\|_2^2}{D}.$$

Now consider the $m = 1$ case where the input belongs to training data. Without loss of generality, assume that the input is the first row of \mathbf{X}_p , i.e. $\mathbf{x}_0 := \mathbf{x}_1$. Note that in this case for $\eta = \mathbf{X}_{\bar{p}} \beta_{\bar{p}} + \epsilon$, we

have $\eta_i \sim \mathcal{N}\left(0, \left(\sigma^2 + \frac{\|\mathbf{x}_{1,p}\|_2^2}{D}\right) \mathbf{I}_n\right)$, for $i = 1$, $\eta_i \sim \mathcal{N}\left(0, \left(1 + \sigma^2 - \frac{p}{D}\right)\right)$, for $i = 2, 3, \dots, n$, and η_i 's are independent. We have

$$\begin{aligned}\widehat{y}_0 &= \mathbf{x}_{1,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p\right)^{-1} \mathbf{X}_p^\top (\mathbf{X}_p \beta_p + \eta) \\ &= \mathbf{x}_{1,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p\right)^{-1} \mathbf{X}_p^\top \mathbf{X}_p \beta_p + \mathbf{x}_{1,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p\right)^{-1} \mathbf{X}_p^\top \eta \\ &= \mathbf{x}_{1,p}^\top \left(\widehat{\Sigma} + \lambda \mathbf{I}_p\right)^{-1} \widehat{\Sigma} \beta_p + \mathbf{x}_{1,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p\right)^{-1} \mathbf{X}_p^\top \eta.\end{aligned}\quad (9)$$

The first term in (9) can be written as

$$\begin{aligned}\mathbf{x}_{1,p}^\top \left(\widehat{\Sigma} + \lambda \mathbf{I}_p\right)^{-1} \widehat{\Sigma} \beta_p &= \mathbf{x}_{1,p}^\top \beta_p - \lambda \mathbf{x}_{1,p}^\top \left(\widehat{\Sigma} + \lambda \mathbf{I}_p\right)^{-1} \beta_p \\ &= \mathbf{x}_{1,p}^\top \beta_p - \lambda \mathbf{x}_{1,p}^\top \left[\frac{1}{n} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top + \mathbf{A}_\lambda\right]^{-1} \beta_p \\ &= \mathbf{x}_{1,p}^\top \beta_p - \lambda \mathbf{x}_{1,p}^\top \left[\mathbf{A}_\lambda^{-1} - \frac{\mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}\right] \beta_p \\ &= \mathbf{x}_{1,p}^\top (\mathbf{I}_p - \lambda \mathbf{A}_\lambda^{-1}) \beta_p + \frac{\lambda \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \beta_p}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \\ &= \mathbf{x}_{1,p}^\top \left(\mathbf{I}_p - \frac{\lambda}{1 + (1/n) \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \mathbf{A}_\lambda^{-1}\right) \beta_p\end{aligned}$$

Hence, since $\beta_i \sim (0, \frac{1}{D} \mathbf{I}_D)$, using Lemma B.2, The first term in (9) converges to a gaussian with zero mean and variance

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{1}{D} \left\| \widehat{\Sigma} \left(\widehat{\Sigma} + \lambda \mathbf{I}_p\right)^{-1} \mathbf{x}_{1,p} \right\|_2^2 &= \lim_{n \rightarrow \infty} \frac{1}{D} \left\| \left(\mathbf{I}_p - \frac{\lambda}{1 + (1/n) \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \mathbf{A}_\lambda^{-1}\right) \mathbf{x}_{1,p} \right\|_2^2 \\ &= \frac{\|\mathbf{x}_{1,p}\|_2^2}{D} \left(1 - \frac{2\lambda g(-\lambda)}{1 + \frac{\gamma \|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda)} + \frac{\lambda^2 g'(-\lambda)}{\left(1 + \frac{\gamma \|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda)\right)^2}\right).\end{aligned}$$

where for the second equality we have used the fact that using Lemma B.1 by setting $C_n = (1/n) \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top$, as $n \rightarrow \infty$, such that $p/n = \gamma$, almost surely,

$$\frac{1}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \rightarrow \frac{1}{n} \|\mathbf{x}_{1,p}\|_2^2 g(-\lambda), \quad \frac{1}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-2} \mathbf{x}_{1,p} \rightarrow \frac{1}{n} \|\mathbf{x}_{1,p}\|_2^2 g'(-\lambda).$$

Now consider the second term in (9). It can be written as

$$\begin{aligned}\mathbf{x}_{1,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p\right)^{-1} \mathbf{X}_p^\top \eta &= \mathbf{x}_{1,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p\right)^{-1} \mathbf{x}_{1,p}^\top \eta_1 \\ &\quad + \sum_{i=2}^n \mathbf{x}_{1,p}^\top \left(\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p\right)^{-1} \mathbf{x}_{i,p}^\top \eta_i \\ &= A \eta_1 + \sum_{i=2}^n B_i \eta_i.\end{aligned}$$

First consider

$$\begin{aligned}
A &= \frac{1}{n} \mathbf{x}_{1,p}^\top \left(\frac{1}{n} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top + \lambda \mathbf{I}_p + \underbrace{\frac{1}{n} \sum_{i=2}^n \mathbf{x}_{i,p} \mathbf{x}_{i,p}^\top}_{\mathbf{A}_\lambda} \right)^{-1} \mathbf{x}_{i,p} \\
&= \frac{1}{n} \mathbf{x}_{1,p}^\top \left(\mathbf{A}_\lambda^{-1} - \frac{\mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \right) \mathbf{x}_{1,p} \\
&= \frac{1}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \left(1 - \frac{\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \right) \\
&= \frac{\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \xrightarrow{a.s.} \frac{\gamma g(-\lambda) (\|\mathbf{x}_{1,p}\|_2^2 / p)}{1 + \gamma g(-\lambda) (\|\mathbf{x}_{1,p}\|_2^2 / p)}.
\end{aligned}$$

Now, consider

$$\begin{aligned}
B_2 &= \frac{1}{n} \mathbf{x}_{1,p}^\top \left(\frac{1}{n} \mathbf{U} \mathbf{U}^\top + \lambda \mathbf{I}_p + \underbrace{\frac{1}{n} \sum_{i=3}^n \mathbf{x}_{i,p} \mathbf{x}_{i,p}^\top}_{\mathbf{A}_{2,\lambda}} \right)^{-1} \mathbf{x}_{i,p}; \quad \mathbf{U} \triangleq [\mathbf{x}_{1,p} \quad \mathbf{x}_{2,p}] \\
&= \frac{1}{n} \mathbf{x}_{1,p}^\top \left[\mathbf{A}_{2,\lambda}^{-1} - \underbrace{\mathbf{A}_{2,\lambda}^{-1} \mathbf{U} \left(n \lambda \mathbf{I}_p + \mathbf{U}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{U} \right)^{-1} \mathbf{U}^\top \mathbf{A}_{2,\lambda}^{-1}}_{\mathbf{C}_2} \right] \mathbf{x}_{2,p}.
\end{aligned}$$

We have

$$\begin{aligned}
\mathbf{C}_2 &= [\mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \quad \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}] \begin{bmatrix} n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} & \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \\ \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} & n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \\ \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \end{bmatrix} \\
&= \left[\underbrace{\left((n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p}) (n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}) - \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right)}_{D_2} \right]^{-1} \tilde{\mathbf{C}}_2.
\end{aligned}$$

We have

$$\begin{aligned}
\tilde{\mathbf{C}}_2 &= [\mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \quad \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}] \begin{bmatrix} n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} & -\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \\ -\mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} & n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \end{bmatrix} \begin{bmatrix} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \\ \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \end{bmatrix} \\
&= \left(n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \right) \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} - \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \\
&\quad - \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} + \left(n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right) \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1}.
\end{aligned}$$

Thus,

$$\begin{aligned}
B_2 &= \frac{1}{n} \mathbf{x}_{1,p}^\top \left[\mathbf{A}_{2,\lambda}^{-1} - \frac{\tilde{\mathbf{C}}_2}{D_2} \right] \mathbf{x}_{2,p} \\
&= \frac{\mathbf{x}_{1,p}^\top}{n} \left\{ \mathbf{A}_{2,\lambda}^{-1} - \left[\left(n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right) \left(n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \right) - \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right]^{-1} \right. \\
&\quad \left[\left(n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \right) \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} - \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \right. \\
&\quad \left. \left. - \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} + \left(n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right) \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \right] \right\} \mathbf{x}_{2,p}.
\end{aligned}$$

using Lemma [B.1](#) by setting $C_n = (1/n)\mathbf{x}\mathbf{x}^\top$, as $n, p \rightarrow \infty$, such that $p/n = \gamma$, almost surely,

$$\frac{1}{n}\mathbf{x}^\top \mathbf{A}_\lambda^{-2}\mathbf{x} \rightarrow \frac{\|\mathbf{x}\|_2^2}{n}g(-\lambda).$$

Hence, letting $n \rightarrow \infty$, B_2 converges weakly to

$$\begin{aligned} B_2' &= \frac{1}{n}\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \left\{ \mathbf{A}_{2,\lambda} - \left[\left(\lambda + g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{n} \right) \left(\lambda + g(-\lambda) \frac{\|\mathbf{x}_{2,p}\|_2^2}{n} \right) - \left(\frac{\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}}{n} \right)^2 \right]^{-1} \right. \\ &\quad \left[\left(\lambda + g(-\lambda) \frac{\|\mathbf{x}_{2,p}\|_2^2}{n} \right) \frac{\mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top}{n} - \left(\frac{\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}}{n} \right) \left(\frac{\mathbf{x}_{1,p} \mathbf{x}_{2,p}^\top}{n} + \frac{\mathbf{x}_{2,p} \mathbf{x}_{1,p}^\top}{n} \right) \right. \\ &\quad \left. \left. + \left(\lambda + g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{n} \right) \frac{\mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top}{n} \right] \right\} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \\ &= \frac{1}{n}\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \frac{\lambda^2}{\left(\lambda + g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{n} \right) \left(\lambda + g(-\lambda) \frac{\|\mathbf{x}_{2,p}\|_2^2}{n} \right) - \left(\frac{\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}}{n} \right)^2}. \end{aligned}$$

Note that by LLN, $(1/n)\|\mathbf{x}_{2,p}\|_2^2 \rightarrow \gamma$ and $(1/n^2)\left(\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}\right)^2 \rightarrow 0$. Further, since $\mathbf{x}_{2,p} \sim (0, \mathbf{I}_p)$, using Lemma [B.2](#) as $n \rightarrow \infty$, $(1/\sqrt{n})\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}$ converges to a gaussian with mean zero and variance

$$\lim_{n \rightarrow \infty} \frac{1}{n} \|\mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p}\|_2^2 = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-2} \mathbf{x}_{1,p} = \frac{1}{n} \|\mathbf{x}_{1,p}\|_2^2 g'(-\lambda)$$

where for the second equality we have used Lemma [B.1](#) with $C_n = (1/n)\mathbf{x}_{1,p}\mathbf{x}_{1,p}^\top$. Hence, B_2 converges weakly to

$$\frac{\tilde{\sigma}}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \frac{\gamma \tilde{\sigma}^2 \|\mathbf{x}_{1,p}\|_2^2}{pn} g'(-\lambda)\right)$$

where

$$\tilde{\sigma} = \frac{\lambda^2}{\lambda^2 + \gamma \lambda g(-\lambda) \left[1 + \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} \right] + \gamma^2 m^2(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{p}}.$$

By symmetry over i , $\sum_{i=2}^n B_i^2 = (n-1)B_2^2$ that converges almost surely to $\frac{\gamma \tilde{\sigma}^2 \|\mathbf{x}_{1,p}\|_2^2}{p} g'(-\lambda)$. Therefore using Lemma [B.2](#)

$$\sum_{i=2}^n B_i \eta_i \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \left(\sigma^2 + 1 - \frac{p}{D}\right) \gamma \tilde{\sigma}^2 \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} g'(-\lambda)\right).$$

Thus, by independence of β_p and η_i 's, as $n, p \rightarrow \infty$ the output y converges in distribution to a gaussian with zero mean and variance

$$\begin{aligned} &\left[\left(\frac{\lambda^2}{(\lambda + \gamma g(-\lambda)) \left(\lambda + \gamma \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda) \right)} \right)^2 \gamma g'(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} \right] \left(\sigma^2 + 1 - \frac{p}{D} \right) \\ &+ \left(\frac{\gamma g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{p}}{1 + \gamma g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{p}} \right)^2 \left(\sigma^2 + \frac{\|\mathbf{x}_{1,p}\|_2^2}{D} \right) \\ &+ \left(1 - \frac{2\lambda g(-\lambda)}{1 + \gamma \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda)} + \frac{\lambda^2 g'(-\lambda)}{\left(1 + \gamma \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda) \right)^2} \right) \frac{\|\mathbf{x}_{1,p}\|_2^2}{D}, \end{aligned}$$

which completes the proof. \square

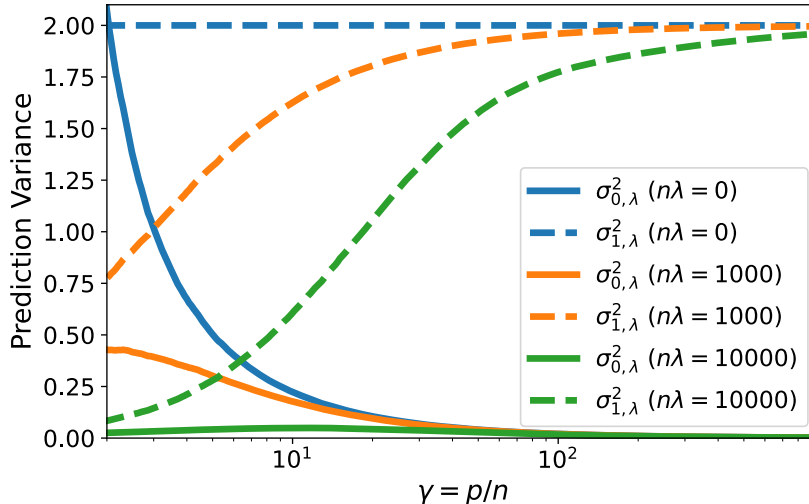


Figure 5: Theoretical variances of the predictions \hat{y}_0 by ridge regularized linear regression models for the Gaussian data setting with $n = 10^3$, $D = 10^7$, and $\sigma = 1$ on a single sampled x_0 for when x_0 is a test point ($\sigma_{0,\lambda}^2$) and when x_0 is a training point ($\sigma_{1,\lambda}^2$) for different amounts of regularization λ . While increased ridge regularization decreases the variance $\sigma_{0,\lambda}^2$ on training point predictions, it also decreases the variance $\sigma_{1,\lambda}^2$ for test points in such a way that the two distributions become easier to distinguish. As such, membership inference is easier for ridge regularized models in this setting.

E Experimental Implementation Details

All experiments were ran only on CPUs on our internal servers without GPU processing. Processors used may have included Intel Xeon CPU E5-2630 (256GB RAM), Intel Xeon Silver 4214 CPU (192GB RAM), Intel Xeon Platinum 8260 CPU (192GB RAM), and AMD Ryzen Threadripper 1900X (32GB RAM). Our code is primarily written in Python and mainly uses numpy implementations of linear algebra operations. Please refer to our code on the Github page for more details.

The histograms in Figure 1 are obtained as follows. We first sample a vector $x_0 \sim \mathcal{N}(0, \mathbf{I}_D)$, where $D = 20,000$. Then, for each $p = \gamma n$, we perform the following procedure 20,000 times. We sample $\beta \sim \mathcal{N}(0, \frac{1}{D}\mathbf{I}_D)$. Then, we sample an $n \times D$ matrix \mathbf{X} such that each element is iid standard normal. We then generate the ground truth vector $\mathbf{y} = \mathbf{X}\beta + \epsilon$, where ϵ is an n -dimensional vector whose elements are iid standard normal. We obtain the least squares estimates $\hat{\beta}$ on the first p columns of \mathbf{X} and on the vector \mathbf{y} using numpy’s ltsq function. Finally, we collect the $\hat{y}_0 = x_{0,p}^\top \hat{\beta}$ of all 20,000 models to form the blue histograms in Figure 1. The orange histograms are formed the same way except that the first row of \mathbf{X} is replaced with x_0 and the first element of \mathbf{y} is replaced with $y_0 = x_0^\top \beta + \epsilon_0$ for $\epsilon_0 \sim \mathcal{N}(0, 1)$.

The experiment in Figure 2b is performed as follows. In the experiment, we estimate the optimal membership advantage. Since the optimal MI adversary requires knowledge of the linear regression model’s output distributions when a data point x_0 is in its training dataset ($m = 1$) and when x_0 is not ($m = 0$), we approximate these distributions by forming discrete histograms. To obtain the samples for the histograms, we use the same procedure as detailed in the previous paragraph, except that we obtain 100,000 samples for each histogram for increased precision. From these samples, the discrete histograms for $(\hat{y}_0 | m = 0)$ and $(\hat{y}_0 | m = 1)$ for a given γ are then formed by splitting the interval between the minimum and maximum values over both $(\hat{y}_0 | m = 0)$ and $(\hat{y}_0 | m = 1)$ into 150 equally spaced bins. The histograms are normalized so that they represent probability mass functions (i.e. the bin counts sum to 1). Finally, treating the two histograms as probability mass functions, the membership advantage is calculated according to Definition 2.1. For Figure 2b, this procedure is repeated 20 times, each with a newly sampled x_0 , and the mean membership advantage over the 20 experiments is plotted.

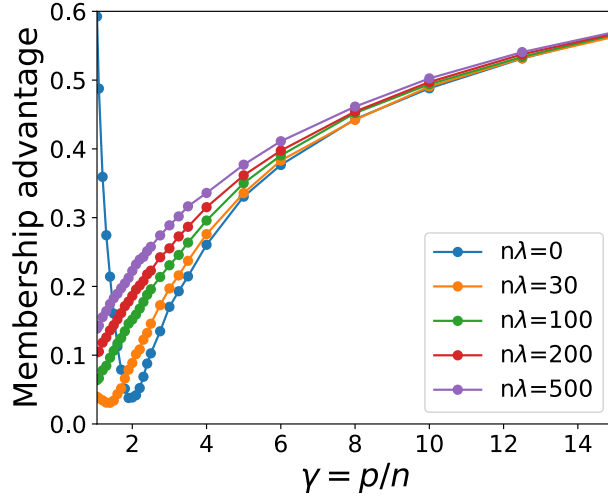


Figure 6: Experimental membership advantages for ridge-regularized linear regression on Gaussian data with $n = 100$, $D = 3000$, and $\sigma = 1$ for different regularization strengths λ . As predicted by our theory, membership advantage increases with additional regularization in the sufficiently overparameterized regime. This experiment verifies our theoretical findings.

The experiments in Figure 4 are also obtained by approximating the optimal MI adversary with discrete histograms as in the previous paragraph. The only difference is in how the datasets (\mathbf{X} , \mathbf{y} , etc.) are sampled. Specifically, they are sampled according to the distributions for each experiment detailed in Section 5. Again, the histograms are formed by splitting the model’s prediction interval for each γ into 150 equally spaced bins. 20 experiments are performed for each data model, with the means and standard errors reported in the figures.

F Experimental Verification of Ridge Theory

We verify our theoretical finding that ridge regression increases membership advantage on linear regression models with Gaussian data in the overparameterized regime. The experiment follows the procedure detailed in Section E for Figure 2b except that we only sample 50,000 datasets for each of $m = 0$ and $m = 1$ for each γ and each λ for computational efficiency. For this experiment, we set $n = 100$, $D = 3,000$, and $\sigma = 1$, as in Figure 2b. The results, shown in Figure 6, closely resemble the trend shown in the theoretical plot in Figure 3a, thus verifying our theory.