

A. Extended Introductory Discussion

Due to space constraint, a few details from the introductory parts (Sections 1,2) were omitted. We bring them in this appendix. We especially recommend the uninformed reader to go over the extended OLS background we provide in Appendix A.3.

A.1. Proof Of Privacy of Algorithm 1

Theorem A.1. *Algorithm 1 is (ϵ, δ) -differentially private.*

Proof. The proof of the theorem is based on the fact the Algorithm 1 is the result of composing the differentially private Propose-Test-Release algorithm of (Dwork & Lei, 2009) with the differentially private analysis of the Johnson-Lindenstrauss transform of (Sheffet, 2015).

More specifically, we use Theorem B.1 from (Sheffet, 2015) that states that given a matrix A whose all of its singular values are greater than $T(\epsilon, \delta)$ where $T(\epsilon, \delta)^2 = \frac{2B^2}{\epsilon} \left(\sqrt{2r \ln(4/\delta)} + 2 \ln(4/\delta) \right)$, publishing RA is (ϵ, δ) -differentially private for a r -row matrix R whose entries sampled are i.i.d normal Gaussians. Since we have that all of the singular values of A' are greater than w (as specified in Algorithm 1), outputting RA' is $(\epsilon/2, \delta/2)$ -differentially private. The rest of the proof boils down to showing that (i) the if-else-condition is $(\epsilon/2, 0)$ -differentially private and that (ii) w.p. $\leq \delta/2$ any matrix A whose smallest singular value is smaller than w passes the if-condition (step 3). If both these facts hold, then knowing whether we pass the if-condition or not is $(\epsilon/2)$ -differentially private and the output of the algorithm is $(\epsilon/2, \delta)$ -differentially private, hence basic composition gives the overall bound of (ϵ, δ) -differential privacy.

To prove (i) we have that for any pair of neighboring matrices A and B that differ only on the i -th row, denoted \mathbf{a}_i and \mathbf{b}_i resp., we have $B^\top B - \mathbf{b}_i \mathbf{b}_i^\top = A^\top A - \mathbf{a}_i \mathbf{a}_i^\top$. Applying Weyl's inequality we have

$$\begin{aligned} \sigma_{\min}(B^\top B) &\leq \sigma_{\min}(B^\top B - \mathbf{b}_i \mathbf{b}_i^\top) + \sigma_{\max}(\mathbf{b}_i \mathbf{b}_i^\top) \\ &\leq \sigma_{\min}(A^\top A) + \sigma_{\max}(\mathbf{a}_i \mathbf{a}_i^\top) + \sigma_{\max}(\mathbf{b}_i \mathbf{b}_i^\top) \\ &\leq \sigma_{\min}(A^\top A) + 2B^2 \end{aligned}$$

hence $|\sigma_{\min}(A)^2 - \sigma_{\min}(B)^2| \leq 2B^2$, so adding $Lap(\frac{4B^2}{\epsilon})$ is $(\epsilon/2)$ -differentially private.

To prove (ii), note that by standard tail-bounds on the Laplace distribution we have that $\Pr[Z < -\frac{4B^2 \ln(1/\delta)}{\epsilon}] \leq \frac{\delta}{2}$. Therefore, w.p. $1 - \delta/2$ it holds that any matrix A that passes the if-test of the algorithm must have $\sigma_{\min}(A)^2 > w^2$. Also note that a similar argument shows that for any $0 < \beta < 1$, any matrix A s.t. $\sigma_{\min}(A)^2 > w^2 + \frac{4B^2 \ln(1/\beta)}{\epsilon}$ passes the if-condition of the algorithm w.p. $1 - \beta$. \square

A.2. Omitted Preliminary Details

Linear Algebra and Pseudo-Inverses. Given a matrix M we denote its SVD as $M = USV^\top$ with U and V being orthonormal matrices and S being a non-negative diagonal matrix whose entries are the singular values of M . We use $\sigma_{\max}(M)$ and $\sigma_{\min}(M)$ to denote the largest and smallest singular value resp. Despite the risk of confusion, we stick to the standard notation of using σ^2 to denote the variance of a Gaussian, and use $\sigma_j(M)$ to denote the j -th singular value of M . We use M^+ to denote the Moore-Penrose inverse of M , defined as $M^+ = VS^{-1}U^\top$ where S^{-1} is a matrix with $S_{j,j}^{-1} = 1/S_{j,j}$ for any j s.t. $S_{j,j} > 0$.

The Gaussian Distribution. A univariate Gaussian $\mathcal{N}(\mu, \sigma^2)$ denotes the Gaussian distribution whose mean is μ and variance σ^2 , with $\text{PDF}(x) = (\sqrt{2\pi\sigma^2})^{-1} \exp(-\frac{x-\mu}{2\sigma^2})$. Standard concentration bounds on Gaussians give that $\Pr[x > \mu + 2\sigma\sqrt{\ln(1/\nu)}] < \nu$ for any $\nu \in (0, \frac{1}{e})$. A multivariate Gaussian $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$ for some positive semi-definite Σ denotes the multivariate Gaussian distribution where the mean of the j -th coordinate is the μ_j and the co-variance between coordinates j and k is $\Sigma_{j,k}$. The PDF of such Gaussian is defined only on the subspace $\text{colspan}(\Sigma)$, where for every $x \in \text{colspan}(\Sigma)$ we have $\text{PDF}(\mathbf{x}) = \left((2\pi)^{\text{rank}(\Sigma)} \cdot \tilde{\det}(\Sigma) \right)^{-1/2} \exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^+ (\mathbf{x} - \boldsymbol{\mu}))$ and $\tilde{\det}(\Sigma)$ is the multiplication of all non-zero singular values of Σ . A matrix Gaussian distribution denoted $\mathcal{N}(M_{a \times b}, U, V)$ has mean M , variance U on its rows and variance V on its columns. For full rank U and V it holds that $\text{PDF}_{\mathcal{N}(M, U, V)}(X) = \frac{(2\pi)^{-ab/2} (\det(U))^{-b/2} (\det(V))^{-a/2}}{\exp(-\frac{1}{2} \text{trace}(V^{-1}(X - M)^\top U^{-1}(X - M))})}$. In our case, we will only use matrix Gaussian distributions with $\mathcal{N}(M_{a \times b}, I_{a \times a}, V)$ and so each row in this matrix is an i.i.d sample from a b -dimensional multivariate Gaussian $\mathcal{N}((M)_{j \rightarrow}, V)$.

We will repeatedly use the rules regarding linear operations on Gaussians. That in, for any c , it holds that $c\mathcal{N}(\mu, \sigma^2) = \mathcal{N}(c \cdot \mu, c^2 \sigma^2)$. For any C it holds that $C \cdot \mathcal{N}(\boldsymbol{\mu}, \Sigma) = \mathcal{N}(C\boldsymbol{\mu}, C\Sigma C^\top)$. And for any C it holds that $\mathcal{N}(M, U, V) \cdot C = \mathcal{N}(MC, U, C^\top V C)$. In particular, for any \mathbf{c} (which can be viewed as a $b \times 1$ -matrix) it holds that $\mathcal{N}(M, U, V) \cdot \mathbf{c} = \mathcal{N}(M\mathbf{c}, U, \mathbf{c}^\top V \mathbf{c}) = \mathcal{N}(M\mathbf{c}, \mathbf{c}^\top V \mathbf{c} \cdot U)$.

We will also require the following proposition.

Proposition A.2. *Given σ^2, λ^2 s.t. $1 \leq \frac{\sigma^2}{\lambda^2} \leq c^2$ for some constant c , let X and Y be two random Gaussians s.t. $X \sim \mathcal{N}(0, \sigma^2)$ and $Y \sim \mathcal{N}(0, \lambda^2)$. It follows that $\frac{1}{c} \text{PDF}_Y(x) \leq \text{PDF}_X(x) \leq c \text{PDF}_{cY}(x)$ for any x .*

Corollary A.3. *Under the same notation as in Proposition A.2, for any set $S \subset \mathbb{R}$ it holds that $\frac{1}{c} \Pr_{x \leftarrow Y}[x \in S] \leq \Pr_{x \leftarrow X}[x \in S] \leq c \Pr_{x \leftarrow cY}[x \in S] =$*

$c\Pr_{x \leftarrow Y}[x \in S/c]$

Proof. The proof is mere calculation.

$$\begin{aligned} \frac{\text{PDF}_X(x)}{\text{PDF}_{cY}(x)} &= \sqrt{\frac{c^2 \lambda^2}{\sigma^2}} \cdot \frac{\exp(-\frac{x^2}{2\sigma^2})}{\exp(-\frac{x^2}{2c^2 \lambda^2})} \\ &\leq c \cdot \exp\left(\frac{x^2}{2} \left(\frac{1}{c^2 \lambda^2} - \frac{1}{\sigma^2}\right)\right) \leq c \cdot \exp(0) = c \\ \frac{\text{PDF}_X(x)}{\text{PDF}_Y(x)} &= \sqrt{\frac{\lambda^2}{\sigma^2}} \cdot \frac{\exp(-\frac{x^2}{2\sigma^2})}{\exp(-\frac{x^2}{2\lambda^2})} \\ &\geq c^{-1} \exp\left(\frac{x^2}{2} \left(\frac{1}{\lambda^2} - \frac{1}{\sigma^2}\right)\right) \geq \frac{\exp(0)}{c} = c^{-1} \end{aligned}$$

□

The T_k -Distribution. The T_k -distribution, where k is referred to as the degrees of freedom of the distribution, denotes the distribution over the reals created by *independently* sampling $Z \sim \mathcal{N}(0, 1)$ and $\|\zeta\|^2 \sim \chi_k^2$, and taking the quantity $\frac{Z}{\sqrt{\|\zeta\|^2/k}}$. Its PDF is given by

$\text{PDF}_{T_k}(x) \propto \left(1 + \frac{x^2}{k}\right)^{-\frac{k+1}{2}}$. It is a known fact that as k increases, T_k becomes closer and closer to a normal Gaussian. The T -distribution is often used to determine suitable bounds on the rate of converges, as we illustrate in Section A.3. As the T -distribution is heavy-tailed, existing tail bounds on the T -distribution (which are of the form: if $\tau_\nu = C\sqrt{k((1/\nu)^{2/k} - 1)}$ for some constant C then $\int_{\tau_\nu}^\infty \text{PDF}_{T_k}(x) dx < \nu$) are often cumbersome to work with. Therefore, in many cases in practice, it common to assume $\nu = \Theta(1)$ (most commonly, $\nu = 0.05$) and use existing tail-bounds on normal Gaussians.

Differential Privacy facts. It is known (Dwork et al., 2006b) that if ALG outputs a vector in \mathbb{R}^d such that for any A and A' it holds that $\|\text{ALG}(A) - \text{ALG}(A')\|_1 \leq B$, then adding Laplace noise $\text{Lap}(1/\epsilon)$ to each coordinate of the output of ALG(A) satisfies ϵ -differential privacy. Similarly, (2006b) showed that if for any neighboring A and A' it holds that $\|\text{ALG}(A) - \text{ALG}(A')\|_2^2 \leq \Delta^2$ then adding Gaussian noise $\mathcal{N}(0, \Delta^2 \cdot \frac{2 \ln(2/\delta)}{\epsilon^2})$ to each coordinate of the output of ALG(A) satisfies (ϵ, δ) -differential privacy.

Another standard result (Dwork et al., 2006a) gives that the composition of the output of a (ϵ_1, δ_1) -differentially private algorithm with the output of a (ϵ_2, δ_2) -differentially private algorithm results in a $(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)$ -differentially private algorithm.

A.3. Detailed Background on Ordinary Least Squares

For the unfamiliar reader, we give a short description of the model under which OLS operates as well as the confidence bounds one derives using OLS. This is by no means an ex-

haustive account of OLS and we refer the interested reader to (Rao, 1973; Muller & Stewart, 2006).

Given n observations $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ where for all i we have $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \mathbb{R}$, we assume the existence of a p -dimensional vector $\beta \in \mathbb{R}^p$ s.t. the label y_i was derived by $y_i = \beta^\top \mathbf{x}_i + e_i$ where $e_i \sim \mathcal{N}(0, \sigma^2)$ independently (also known as the homoscedastic Gaussian model). We use the matrix notation where X denotes the $(n \times p)$ -matrix whose rows are \mathbf{x}_i , and use $\mathbf{y}, \mathbf{e} \in \mathbb{R}^n$ to denote the vectors whose i -th entry is y_i and e_i resp. To simplify the discussion, we assume X has full rank.

The parameters of the model are therefore β and σ^2 , which we set to discover. To that end, we minimize $\min_{\mathbf{z}} \|\mathbf{y} - X\mathbf{z}\|^2$ and solve

$$\hat{\beta} = (X^\top X)^{-1} X^\top \mathbf{y} = (X^\top X)^{-1} X^\top (X\beta + \mathbf{e}) = \beta + X^+ \mathbf{e}$$

As $\mathbf{e} \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 I_{n \times n})$, it holds that $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^\top X)^{-1})$, or alternatively, that for every coordinate j it holds that $\hat{\beta}_j = \mathbf{e}_j^\top \hat{\beta} \sim \mathcal{N}(\beta_j, \sigma^2 (X^\top X)^{-1}_{j,j})$.

Hence we get $\frac{\hat{\beta}_j - \beta_j}{\sigma \sqrt{(X^\top X)^{-1}_{j,j}}} \sim \mathcal{N}(0, 1)$. In addition, we denote the vector

$$\zeta = \mathbf{y} - X\hat{\beta} = (X\beta + \mathbf{e}) - X(\beta + X^+ \mathbf{e}) = (I - XX^+) \mathbf{e}$$

and since XX^+ is a rank- p (symmetric) projection matrix, we have $\zeta \sim \mathcal{N}(0, \sigma^2 (I - XX^+))$. Therefore, $\|\zeta\|^2$ is equivalent to summing the squares of $(n - p)$ i.i.d samples from $\mathcal{N}(0, \sigma^2)$. In other words, the quantity $\|\zeta\|^2 / \sigma^2$ is sampled from a χ^2 -distribution with $(n - p)$ degrees of freedom.

We sidetrack from the OLS discussion to give the following bounds on the l_2 -distance between β and $\hat{\beta}$, as the next claim shows.

Claim A.4. For any $0 < \nu < 1/2$, the following holds w.p. $\geq 1 - \nu$ over the randomness of the model (the randomness over \mathbf{e})

$$\begin{aligned} \|\beta - \hat{\beta}\|^2 &= \|X^+ \mathbf{e}\|^2 \\ &= O(\sigma^2 \log(p/\nu) \cdot \|X^+\|_F^2) \end{aligned} \quad (6)$$

$$\begin{aligned} \|\hat{\beta}\|^2 &= \|\beta + X^+ \mathbf{e}\|^2 \\ &= O\left(\|\beta\| + \sigma \cdot \|X^+\|_F \cdot \sqrt{\log(p/\nu)}\right)^2 \end{aligned}$$

$$\left| \frac{1}{n-p} \|\zeta\|^2 - \sigma^2 \right| = O\left(\sqrt{\frac{\ln(1/\nu)}{n-p}}\right)$$

Proof. Since $\mathbf{e} \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 I_{n \times n})$ then $X^+ \mathbf{e} \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 (X^\top X)^{-1})$. Denoting the SVD decomposition $(X^\top X)^{-1} = VSV^\top$ with S denoting the diagonal matrix whose entries are $\sigma_{\max}^{-2}(X), \dots, \sigma_{\min}^{-2}(X)$, we have that $V^\top X^+ \mathbf{e} \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 S)$. And so, each coordinate of $V^\top X^+ \mathbf{e}$ is distributed like an i.i.d Gaussian. So

w.p. $\geq 1 - \nu/2$ non of these Gaussians is a factor of $O(\sigma\sqrt{\ln(p/\nu)})$ greater than its standard deviation. And so w.p. $\geq 1 - \nu/2$ it holds that $\|X^+\mathbf{e}\|^2 = \|V^\top X^+\mathbf{e}\|^2 \leq O(\sigma^2 \log(p/\nu) (\sum_i \sigma_i^{-2}(X)))$. Since $\sum_i \sigma_i^{-2}(X) = \text{trace}((X^\top X)^{-1}) = \text{trace}(X^+(X^+)^\top) = \|X^+\|_F^2$, the bound of (6) is proven.

The bound on $\|\hat{\beta}\|^2$ is an immediate corollary of (6) using the triangle inequality.⁸ The bound on $\|\zeta\|^2$ follows from tail bounds on the χ_{n-p}^2 distribution, as detailed in Section 2. \square

Returning to OLS, it is important to note that $\hat{\beta}$ and ζ are independent of one another. (Note, $\hat{\beta}$ depends solely on $X^+\mathbf{e} = (X^+X)X^+\mathbf{e} = X^+P_U\mathbf{e}$, whereas ζ depends on $(I - XX^+)\mathbf{e} = P_{U^\perp}\mathbf{e}$. As \mathbf{e} is spherically symmetric, the two projections are independent of one another and so $\hat{\beta}$ is independent of ζ .) As a result of the above two calculations, we have that the quantity

$$t_{\hat{\beta}_j}(\beta_j) \stackrel{\text{def}}{=} \frac{\hat{\beta}_j - \beta_j}{\sqrt{(X^\top X)_{j,j}^{-1} \cdot \frac{\|\zeta\|^2}{n-p}}} = \frac{\hat{\beta}_j - \beta_j}{\sigma \sqrt{(X^\top X)_{j,j}^{-1}}} \bigg/ \frac{\|\zeta\|}{\sigma \sqrt{n-p}}$$

is distributed like a T -distribution with $(n-p)$ degrees of freedom. Therefore, we can compute an exact probability estimation for this quantity. That is, for any measurable $S \subset \mathbb{R}$ we have

$$\Pr \left[\hat{\beta} \text{ and } \zeta \text{ satisfying } t_{\hat{\beta}_j}(\beta_j) \in S \right] = \int_S \text{PDF}_{T_{n-p}}(x) dx$$

The importance of the t -value $t(\beta_j)$ lies in the fact that it can be fully estimated from the observed data X and \mathbf{y} (for any value of β_j), which makes it a *pivotal quantity*. Therefore, given X and \mathbf{y} , we can use $t(\beta_j)$ to describe the likelihood of any β_j — for any $z \in \mathbb{R}$ we can now give an estimation of how likely it is to have $\beta_j = z$ (which is $\text{PDF}_{T_{n-p}}(t(z))$). The t -values enable us to perform multitude of statistical inferences. For example, we can say which of two hypotheses is more likely and by how much (e.g., we are 5-times more likely that the hypothesis $\beta_j = 3$ is true than the hypothesis $\beta_j = 14$ is true); we can compare between two coordinates j and j' and report we are more confident that $\beta_j > 0$ than $\beta_{j'} > 0$; or even compare among the t -values we get across multiple datasets (such as the datasets we get from subsampling rows from a single dataset).

In particular, we can use $t(\beta_j)$ to α -reject unlikely values of β_j . Given $0 < \alpha < 1$, we denote c_α as the number for which the interval $(-c_\alpha, c_\alpha)$ contains a probability mass of $1 - \alpha$ from the T_{n-p} -distribution. And so we derive a

⁸Observe, though \mathbf{e} is spherically symmetric, and is likely to be approximately-orthogonal to β , this does not necessarily hold for $X^+\mathbf{e}$ which isn't spherically symmetric. Therefore, we result to bounding the l_2 -norm of $\hat{\beta}$ using the triangle bound.

corresponding *confidence interval* I_α centered at $\hat{\beta}_j$ where $\beta_j \in I_\alpha$ with confidence of level of $1 - \alpha$.

We comment as to the actual meaning of this confidence interval. Our analysis thus far applied w.h.p to a vector \mathbf{y} derived according to this model. Such X and \mathbf{y} will result in the quantity $t_{\hat{\beta}_j}(\beta_j)$ being distributed like a T_{n-p} -distribution — where β_j is given as the model parameters and $\hat{\beta}_j$ is the random variable. We therefore have that guarantee that for X and \mathbf{y} derived according to this model, the event $E_\alpha \stackrel{\text{def}}{=} \hat{\beta}_j \in \left(\beta_j \pm c_\alpha \cdot \sqrt{(X^\top X)_{j,j}^{-1} \cdot \frac{\|\zeta\|^2}{n-p}} \right)$ happens w.p. $1 - \alpha$. However, the analysis done over a *given* dataset X and \mathbf{y} (once \mathbf{y} has been drawn) views the quantity $t_{\hat{\beta}_j}(\beta_j)$ with $\hat{\beta}_j$ given and β_j unknown. Therefore the event E_α either holds or does not hold. That is why the alternative terms of *likelihood* or *confidence* are used, instead of probability. We have a confidence level of $1 - \alpha$ that indeed $\beta_j \in \hat{\beta}_j \pm c_\alpha \cdot \sqrt{(X^\top X)_{j,j}^{-1} \cdot \frac{\|\zeta\|^2}{n-p}}$, because this event does happen in $1 - \alpha$ fraction of all datasets generated according to our model.

Rejecting the Null Hypothesis. One important implication of the quantity $t(\beta_j)$ is that we can refer specifically to the hypothesis that $\beta_j = 0$, called the *null hypothesis*. This quantity, $t_0 \stackrel{\text{def}}{=} t_{\hat{\beta}_j}(0) = \frac{\hat{\beta}_j \sqrt{n-p}}{\|\zeta\| \sqrt{(X^\top X)_{j,j}^{-1}}}$, represents how

large is $\hat{\beta}_j$ relatively to the empirical estimation of standard deviation σ . Since it is known that as the number of degrees of freedom of a T -distribution tends to infinity then the T -distribution becomes a normal Gaussian, it is common to think of t_0 as a sample from a normal Gaussian $\mathcal{N}(0, 1)$. This allows us to associate t_0 with a p -value, estimating the event “ β_j and $\hat{\beta}_j$ have different signs.” Formally, we define $p_0 = \int_{|t_0|}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$. It is common to reject the null hypothesis when p_0 is sufficiently small (typically, below 0.05).⁹

Specifically, given $\alpha \in (0, 1/2)$, we say we α -reject the *null hypothesis* if $p_0 < \alpha$. Let τ_α be the number s.t. $\Phi(\tau_\alpha) = \int_{\tau_\alpha}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \alpha$. (Standard bounds give that $\tau_\alpha < 2\sqrt{\ln(1/\alpha)}$.) This means we α -reject the null hypothesis if $t_0 > \tau_\alpha$ or $t_0 < -\tau_\alpha$, meaning if $|\hat{\beta}_j| > \tau_\alpha \sqrt{(X^\top X)_{j,j}^{-1} \cdot \frac{\|\zeta\|^2}{n-p}}$.

We can now lower bound the number of i.i.d sample points needed in order to α -reject the null hypothesis. This bound will be our basis for comparison — between standard OLS and the differentially private version.¹⁰

⁹Indeed, it is more accurate to associate with t_0 the value $\int_{|t_0|}^{\infty} \text{PDF}_{T_{n-p}}(x) dx$ and check that this value is $< \alpha$. However, as most uses take α to be a constant (often $\alpha = 0.05$), asymptotically the threshold we get for rejecting the null hypothesis are the same.

¹⁰This theorem is far from being new (except for maybe fo-

Theorem A.5 (Theorem 2.2 restated.). *Fix any positive definite matrix $\Sigma \in \mathbb{R}^{p \times p}$ and any $\nu \in (0, \frac{1}{2})$. Fix parameters $\beta \in \mathbb{R}^p$ and σ^2 and a coordinate j s.t. $\beta_j \neq 0$. Let X be a matrix whose n rows are i.i.d samples from $\mathcal{N}(\mathbf{0}, \Sigma)$, and \mathbf{y} be a vector where $y_i - (X\beta)_i$ is sampled i.i.d from $\mathcal{N}(0, \sigma^2)$. Fix $\alpha \in (0, 1)$. Then w.p. $\geq 1 - \nu$ we have that the $(1 - \alpha)$ -confidence interval is of length $O(c_\alpha \sqrt{\sigma^2 / (n\sigma_{\min}(\Sigma))})$ provided $n \geq C_1(p + \ln(1/\nu))$ for some sufficiently large constant C_1 . Furthermore, there exists a constant C_2 such that w.p. $\geq 1 - \alpha - \nu$ we (correctly) reject the null hypothesis provided*

$$n \geq \max \left\{ C_1(p + \ln(1/\nu)), C_2 \frac{\sigma^2}{\beta_j^2} \cdot \frac{c_\alpha^2 + \tau_\alpha^2}{\sigma_{\min}(\Sigma)} \right\}$$

Here c_α denotes the number for which $\int_{-c_\alpha}^{c_\alpha} \text{PDF}_{T_{n-p}}(x) dx = 1 - \alpha$. (If we are content with approximating T_{n-p} with a normal Gaussian than one can set $c_\alpha \approx \tau_\alpha < 2\sqrt{\ln(1/\alpha)}$.)

Proof. The discussion above shows that w.p. $\geq 1 - \alpha$ we have $|\beta_j - \hat{\beta}_j| \leq c_\alpha \sqrt{(X^\top X)_{j,j}^{-1} \frac{\|\zeta\|^2}{n-p}}$; and in order to α -reject the null hypothesis we must have $|\hat{\beta}_j| > \tau_\alpha \sqrt{(X^\top X)_{j,j}^{-1} \frac{\|\zeta\|^2}{n-p}}$. Therefore, a sufficient condition for OLS to α -reject the null-hypothesis is to have n large enough s.t. $|\beta_j| > (c_\alpha + \tau_\alpha) \sqrt{(X^\top X)_{j,j}^{-1} \frac{\|\zeta\|^2}{n-p}}$. We therefore argue that w.p. $\geq 1 - \nu$ this inequality indeed holds.

We assume each row of X i.i.d vector $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}_p, \Sigma)$, and recall that according to the model $\|\zeta\|^2 \sim \sigma^2 \chi^2(n-p)$. Straightforward concentration bounds on Gaussians and on the χ^2 -distribution give:

(i) W.p. $\leq \alpha$ it holds that $\|\zeta\| > \sigma(\sqrt{n-p} + 2\ln(2/\alpha))$. (This is part of the standard OLS analysis.)

(ii) W.p. $\leq \nu$ it holds that $\sigma_{\min}(X^\top X) \leq \sigma_{\min}(\Sigma)(\sqrt{n} - (\sqrt{p} + \sqrt{2\ln(2/\nu)}))^2$. (Rudelson & Vershynin, 2009)

Therefore, due to the lower bound $n = \Omega(p + \ln(1/\nu))$, w.p. $\geq 1 - \nu - \alpha$ we have that none of these events hold. In such a case we have $\sqrt{(X^\top X)_{j,j}^{-1}} \leq \sqrt{\sigma_{\max}((X^\top X)^{-1})} = O(\frac{1}{\sqrt{n\sigma_{\min}(\Sigma)}})$ and $\|\zeta\| = O(\sigma\sqrt{n-p})$. This implies that the confidence interval of level $1 - \alpha$ has length of $c_\alpha \sqrt{(X^\top X)_{j,j}^{-1} \cdot \frac{\|\zeta\|^2}{n-p}} = O\left(c_\alpha \sqrt{\frac{\sigma^2}{n\sigma_{\min}(\Sigma)}}\right)$; and that in order to α -reject that null-hypothesis it suffices to have $|\beta_j| = \Omega\left((c_\alpha + \tau_\alpha) \sqrt{\frac{\sigma^2}{n\sigma_{\min}(\Sigma)}}\right)$. Plugging in the lower bound on n , we see that this inequality holds.

We comment that for sufficiently large constants C_1, C_2 ,

cusing on the setting where every row in X is sampled from an i.i.d multivariate Gaussians), it is just stated in a non-standard way, discussing solely the power of the t -test in OLS. For further discussions on sample size calculations see (Muller & Stewart, 2006).

it holds that all the constants hidden in the O - and Ω -notations of the proof are close to 1. I.e., they are all within the interval $(1 \pm \eta)$ for some small $\eta > 0$ given $C_1, C_2 \in \Omega(\eta^{-2})$. \square

B. Projecting the Data using Gaussian Johnson-Lindenstrauss Transform

B.1. Main Theorem Restated and Further Discussion

Theorem B.1 (Theorem 3.1 restated.). *Let X be a $n \times p$ matrix, and parameters $\beta \in \mathbb{R}^p$ and σ^2 are such that we generate the vector $\mathbf{y} = X\beta + \mathbf{e}$ with each coordinate of \mathbf{e} sampled independently from $\mathcal{N}(0, \sigma^2)$. Assume $\sigma_{\min}(X) \geq C \cdot w$ and that n is sufficiently large s.t. all of the singular values of the matrix $[X; \mathbf{y}]$ are greater than $C \cdot w$ for some large constant C , and so Algorithm 1 projects the matrix $A = [X; \mathbf{y}]$ without altering it, and publishes $[RX; R\mathbf{y}]$.*

Fix $\nu \in (0, 1/2)$ and $r = p + \Omega(\ln(1/\nu))$. Fix coordinate j . Then w.p. $\geq 1 - \nu$ we have that deriving $\tilde{\beta}$, $\tilde{\zeta}$ and $\tilde{\sigma}^2$ as follows

$$\begin{aligned} \tilde{\beta} &= (X^\top R^\top R X)^{-1} (R X)^\top (R \mathbf{y}) = \beta + (R X)^+ \text{Re} \\ \tilde{\zeta} &= \frac{1}{\sqrt{r}} R \mathbf{y} - \frac{1}{\sqrt{r}} (R X) \tilde{\beta} \\ &= \frac{1}{\sqrt{r}} (I - (R X) (X^\top R^\top R X)^{-1} (R X)^\top) \text{Re} \\ \tilde{\sigma}^2 &= \frac{r}{r-p} \|\tilde{\zeta}\|^2 \end{aligned}$$

then the pivot quantity

$$\tilde{t}(\beta_j) = \frac{\tilde{\beta}_j - \beta_j}{\tilde{\sigma} \sqrt{(X^\top R^\top R X)_{j,j}^{-1}}}$$

has a distribution \mathcal{D} satisfying $e^{-a} \text{PDF}_{T_{r-p}}(x) \leq \text{PDF}_{\mathcal{D}}(x) \leq e^a \text{PDF}_{T_{r-p}}(e^{-a}x)$ for any $x \in \mathbb{R}$, where we denote $a = \frac{r-p}{n-p}$.

Comparison with Existing Bounds. Sarlos' work (2006) utilizes the fact that when r , the numbers of rows in R , is large enough, then $\frac{1}{\sqrt{r}}R$ is a Johnson-Lindenstrauss matrix. Specifically, given r and $\nu \in (0, 1)$ we denote $\eta = \Omega(\sqrt{\frac{p \ln(p) \ln(1/\nu)}{r}})$, and so $r = O(\frac{p \ln(p) \ln(1/\nu)}{\eta^2})$. Let us denote $\tilde{\beta} = \arg \min_{\mathbf{z}} \frac{1}{r} \|R X \mathbf{z} - R \mathbf{y}\|^2$. In this setting, Sarlos' work (Sarlos, 2006) (Theorem 12(3)) guarantees that w.p. $\geq 1 - \nu$ we have $\|\hat{\beta} - \tilde{\beta}\|_2 \leq \eta \|\zeta\| / \sigma_{\min}(X) = O\left(\sqrt{\frac{p \log(p) \log(1/\nu)}{r \sigma_{\min}(X^\top X)}} \|\zeta\|\right)$. Naively bounding $|\hat{\beta}_j - \tilde{\beta}_j| \leq \|\hat{\beta} - \tilde{\beta}\|$ and using the confidence interval for $\hat{\beta}_j - \beta_j$ from Section A.3¹¹

¹¹Where we approximate c_α , the tail bound of the T_{n-p} -distribution with the tail bound on a Gaussian, i.e., use the approximation $c_\alpha \approx O(\sqrt{\ln(1/\alpha)})$.

gives a confidence interval of level $1 - (\alpha + \nu)$ centered at $\tilde{\beta}_j$ with length of $O\left(\sqrt{\frac{p \ln(p) \log(1/\nu)}{r \sigma_{\min}(X^T X)}} \|\zeta\|\right) + O\left(\sqrt{\frac{p \ln(p) \log(1/\alpha)}{r \sigma_{\min}(X^T X)}} \|\zeta\|\right) = O\left(\sqrt{\frac{p \ln(p) \log(1/\nu) + \log(1/\alpha)}{r \sigma_{\min}(X^T X)}} \|\zeta\|\right)$. This implies that our confidence interval has decreased its degrees of freedom from $n - p$ to roughly $r/p \ln(p)$, and furthermore, that it no longer depends on $(X^T X)^{-1}_{j,j}$ but rather on $1/\sigma_{\min}(X^T X)$. It is only due to the fact that we rely on Gaussians and by mimicking carefully the original proof that we can deduce that the t -value has (roughly) $r - p$ degrees of freedom and depends solely on $(X^T X)^{-1}_{j,j}$.

(In the worst case, we have that $(X^T X)^{-1}_{j,j}$ is proportional to $\sigma_{\min}(X^T X)^{-1}$, but it is not uncommon to have matrices where the former is much larger than the latter.) As mentioned in the introduction, alternative techniques ((Chaudhuri et al., 2011; Bassily et al., 2014; Ullman, 2015)) for finding a DP estimator β^{dp} of the linear regression give a data-independent¹² bound of $\|\beta^{dp} - \hat{\beta}\| = \tilde{O}(p/\epsilon)$. Such bounds are harder to compare with the interval length given by Corollary 3.2. Indeed, as we discuss in Section 3 under “Rejecting the null-hypothesis,” enough samples from a multivariate Gaussian whose covariance-matrix is well conditioned give a bound which is well below the worst-upper bound of $O(p/\epsilon)$. (Yet, it is possible that these techniques also do much better on such “well-behaved” data.) What the works of Sarlos and alternative works regrading differentially private linear regression do not take into account are questions such as generating a likelihood for β_j nor do they discuss rejecting the null hypothesis.

B.2. Proof of Theorem 3.1

We now turn to our analysis of $\tilde{\beta}$ and $\tilde{\zeta}$, where our goal is to show that the distribution of the t -values as specified in Theorem 3.1 is well-approximated by the T_{r-p} -distribution. For now, we assume the existence of fixed vectors $\beta \in \mathbb{R}^p$ and $e \in \mathbb{R}^n$ s.t. $y = X\beta + e$. (Later, we will return to the homoscedastic model where each coordinate of e is sampled i.i.d from $\mathcal{N}(0, \sigma^2)$) In other words, we first examine the case where R is the sole source of randomness in our estimation. Based on the assumption that e is fixed, we argue the following.

Claim B.2. *In our model, given X and the output $M = RX$, we have that $\tilde{\beta} \sim \mathcal{N}(\beta + X^+e, \|P_{U^\perp}e\|^2(M^T M)^{-1})$ and $\tilde{\zeta} \sim \mathcal{N}(\mathbf{0}_n, \frac{\|P_{U^\perp}e\|^2}{r}(I_{r \times r} - M(M^T M)^{-1}M^T))$. Where P_{U^\perp} denotes the projection operator onto the subspace orthogonal to $\text{colspan}(X)$; i.e., $P_U = XX^+$ and $P_{U^\perp} = (I_{r \times r} - XX^+)$.*

¹²In other words, independent of X, ζ .

Proof. The matrix R is sampled from $\mathcal{N}(0_{r \times p}, I_{r \times r}, I_{p \times p})$. Given X and $RX = M$, we learn the projection of each row in R onto the subspace spanned by the columns of X . That is, denoting u^T as the i -th row of R and v^T as the i -th row of M , we have that $X^T u = v$. Recall, initially $u \sim \mathcal{N}(\mathbf{0}_n, I_{n \times n})$ – a spherically symmetric Gaussian. As a result, we can denote $u = P_U u + P_{U^\perp} u$ where the two projections are independent samples from $\mathcal{N}(\mathbf{0}_n, P_U)$ and $\mathcal{N}(\mathbf{0}_n, P_{U^\perp})$ resp. However, once we know that $v = X^T u$ we have that $P_U u = X(X^T X)^{-1}X^T u = X(X^T X)^{-1}v$ so we learn $P_U u$ exactly, whereas we get no information about $P_{U^\perp} u$ so $P_{U^\perp} u$ is still sampled from a Gaussian $\mathcal{N}(\mathbf{0}_n, P_{U^\perp})$. As we know for each row of R that $u^T P_U = v^T X^+$, we therefore have that

$$R = RP_U + RP_{U^\perp} = MX^+ + RP_{U^\perp}$$

where $RP_{U^\perp} \sim \mathcal{N}(0_{r \times n}, I_{r \times r}, P_{U^\perp})$. From here on, we just rely on the existing results about the linearity of Gaussians.

$$\begin{aligned} R &\sim \mathcal{N}(MX^+, I_{r \times r}, P_{U^\perp}) \\ &\Rightarrow Re \sim \mathcal{N}(MX^+e, \|P_{U^\perp}e\|^2 I_{r \times r}) \\ &\Rightarrow M^+ Re \sim \mathcal{N}(X^+e, \|P_{U^\perp}e\|^2 (M^T M)^{-1}) \end{aligned}$$

so $\tilde{\beta} = \beta + M^+ Re$ implies $\tilde{\beta} \sim \mathcal{N}(\beta + X^+e, \|P_{U^\perp}e\|^2 (M^T M)^{-1})$. And as $\tilde{\zeta} = \frac{1}{\sqrt{r}}(I_{r \times r} - M(M^T M)^{-1}M^T)Re$ then we have $\tilde{\zeta} \sim \mathcal{N}(\mathbf{0}_r, \frac{\|P_{U^\perp}e\|^2}{r}(I_{r \times r} - MM^+))$ as $(I_{r \times r} - MM^+)M = 0_{r \times p}$. \square

Claim B.2 was based on the assumption that e is fixed. However, given X and y there are many different ways to assign vectors β and e s.t. $y = X\beta + e$. However, the distributions we get in Claim B.2 are *unique*. To see that, recall Equations (1) and (2): $\beta + X^+e = X^+y = \hat{\beta}$ and $P_{U^\perp}e = P_{U^\perp}y = (I - XX^+)y = \zeta$. We therefore have $\tilde{\beta} \sim \mathcal{N}(\hat{\beta}, \|\zeta\|^2 (M^T M)^{-1})$ and $\tilde{\zeta} \sim \mathcal{N}(\mathbf{0}_n, \frac{\|\zeta\|^2}{r}(I - MM^+))$. We will discuss this further, in Section 4, where we will not be able to better analyze the explicit distributions of our estimators. But in this section, we are able to argue more about the distributions of $\tilde{\beta}$ and $\tilde{\zeta}$.

So far we have considered the case that e is fixed, whereas our goal is to argue about the case where each coordinate of e is sampled i.i.d from $\mathcal{N}(0, \sigma^2)$. To that end, we now switch to an intermediate model, in which $P_U e$ is sampled from a multivariate Gaussian while $P_{U^\perp} e$ is fixed as some arbitrary vector of length l . Formally, let \mathcal{D}_l denote the distribution where $P_U e \sim \mathcal{N}(0, \sigma^2 P_U)$ and $P_{U^\perp} e$ is fixed as some specific vector whose length is denoted by $\|P_{U^\perp} e\| = l$.

Claim B.3. *Under the same assumptions as in Claim B.2, given that $e \sim \mathcal{D}_l$, we have that*

$$\begin{aligned} \tilde{\beta} &\sim \mathcal{N}(\beta, \sigma^2(X^T X)^{-1} + l^2(M^T M)^{-1}) \quad \text{and} \\ \tilde{\zeta} &\sim \mathcal{N}\left(\mathbf{0}_n, \frac{l^2}{r}(I - MM^+)\right). \end{aligned}$$

Proof. Recall, $\tilde{\beta} = \beta + X^+ \mathbf{e} + M^+ R(P_{U^\perp} \mathbf{e}) = \beta + M^+(MX^+ + RP_{U^\perp})\mathbf{e} = \beta + X^+ \mathbf{e} + M^+ R(P_{U^\perp} \mathbf{e})$. Now, under the assumption $\mathbf{e} \sim \mathcal{D}_l$ we have that β is the sum of two independent Gaussians:

$$\begin{aligned} \beta + X^+ \mathbf{e} &\sim \mathcal{N}(\beta, \sigma^2(X^+ \cdot P_U \cdot (X^+)^T)) \\ &= \mathcal{N}(\beta, \sigma^2(X^T X)^{-1}) \\ RP_{U^\perp} \mathbf{e} &\sim \mathcal{N}(\mathbf{0}_r, \|P_{U^\perp} \mathbf{e}\|^2 I_{r \times r}) \\ \Rightarrow M^+ R \mathbf{e} &\sim \mathcal{N}(\mathbf{0}_p, \|P_{U^\perp} \mathbf{e}\|^2 (M^T M)^{-1}) \end{aligned}$$

Summing the two independent Gaussians' means and variances gives the distribution of $\tilde{\beta}$. Furthermore, in Claim B.2 we have already established that for any fixed \mathbf{e} we have $\tilde{\zeta} \sim \mathcal{N}\left(\mathbf{0}_n, \frac{\|P_{U^\perp} \mathbf{e}\|^2}{r}(I - MM^+)\right)$. Hence, for $\mathbf{e} \sim \mathcal{D}_l$ we still have $\tilde{\zeta} \sim \mathcal{N}\left(\mathbf{0}_n, \frac{l^2}{r}(I - MM^+)\right)$. (It is easy to verify that the same chain of derivations is applicable when $\mathbf{e} \sim \mathcal{D}_l$.) \square

Corollary B.4. *Given that $\mathbf{e} \sim \mathcal{D}_l$ we have that $\tilde{\beta}_j \sim \mathcal{N}(\beta_j, \sigma^2(X^T X)_{j,j}^{-1} + l^2(M^T M)_{j,j}^{-1})$ for any coordinate j , and that $\|\tilde{\zeta}\|^2 \sim \frac{l^2}{r} \cdot \chi_{r-p}^2$.*

Proof. The corollary follows immediately from the fact that $\beta_j = \mathbf{e}_j^T \tilde{\beta}$, and from the definition of the χ^2 -distribution, as $\tilde{\zeta}$ is a spherically symmetric Gaussian defined on the subspace $\text{colspan}(M)^\perp$ of dimension $r - p$. \square

To continue, we need the following claim.

Claim B.5. *Given X and $M = RX$, and given that $\mathbf{e} \sim \mathcal{D}_l$ we have that $\tilde{\beta}$ and $\tilde{\zeta}$ are independent.*

Proof. Recall, $\tilde{\beta} = \beta + X^+ \mathbf{e} + M^+ R(P_{U^\perp} \mathbf{e})$. And so, given X , M and a specific vector $P_{U^\perp} \mathbf{e}$ we have that the distribution of $\tilde{\beta}$ depends on (i) the projection of \mathbf{e} on $U = \text{colspan}(X)$ and on (ii) the projection of each row in R onto $\tilde{U} = \text{colspan}(M)$. The distribution of $\tilde{\zeta} = \frac{1}{\sqrt{r}} P_{U^\perp} R \mathbf{e} = \frac{1}{\sqrt{r}} P_{U^\perp} (MX^+ + RP_{U^\perp}) \mathbf{e} = \frac{1}{\sqrt{r}} P_{U^\perp} RP_{U^\perp} \mathbf{e}$ depends on (i) the projection of \mathbf{e} onto U^\perp (which for the time being is fix to some specific vector of length l) and on (ii) the projection of each row in R onto \tilde{U}^\perp . Since $P_U \mathbf{e}$ is independent from $P_{U^\perp} \mathbf{e}$, and since for any row \mathbf{u}^T of R we have that $P_{\tilde{U}} \mathbf{u}$ is independent of $P_{\tilde{U}^\perp} \mathbf{u}$, and since \mathbf{e} and R are chosen independently, we have that $\tilde{\beta}$ and $\tilde{\zeta}$ are independent.

Formally, consider any pair of coordinates $\tilde{\beta}_j$ and $\tilde{\zeta}_k$, and we have

$$\tilde{\beta}_j - \beta_j = \mathbf{e}_j^T X^+ \mathbf{e} + \mathbf{e}_j^T M^+ (RP_{U^\perp} \mathbf{e})$$

$$\tilde{\zeta}_k = \mathbf{e}_k^T P_{\tilde{U}^\perp} (RP_{U^\perp} \mathbf{e})$$

Recall, we are given X and $M = RX$. Therefore, we know P_U and $P_{\tilde{U}}$. And so

$$\begin{aligned} \text{Cov}[\tilde{\beta}_j, \tilde{\zeta}_k] &= \mathbf{E}[(\tilde{\beta}_j - \beta_j)(\tilde{\zeta}_k - 0)] \\ &= \mathbf{E}[\mathbf{e}_j^T X^+ \mathbf{e} (RP_{U^\perp} \mathbf{e})^T P_{\tilde{U}^\perp} \mathbf{e}_k] \\ &\quad + \mathbf{E}[\mathbf{e}_j^T M^+ (RP_{U^\perp} \mathbf{e}) (RP_{U^\perp} \mathbf{e})^T P_{\tilde{U}^\perp} \mathbf{e}_k] \\ &= \mathbf{e}_j^T X^+ \mathbf{E}[\mathbf{e} \mathbf{e}^T P_{U^\perp}] \mathbf{E}[R^T] P_{\tilde{U}^\perp} \mathbf{e}_k \\ &\quad + \mathbf{e}_j^T M^+ \mathbf{E}[(RP_{U^\perp} \mathbf{e}) (RP_{U^\perp} \mathbf{e})^T] P_{\tilde{U}^\perp} \mathbf{e}_k \\ &= \mathbf{e}_j^T X^+ \mathbf{E}[\mathbf{e} \mathbf{e}^T P_{U^\perp}] ((MX^+)^T + \mathbf{E}[(RP_{U^\perp})^T]) P_{\tilde{U}^\perp} \mathbf{e}_k \\ &\quad + \mathbf{e}_j^T M^+ (\|P_{U^\perp} \mathbf{e}\|^2 I_{r \times r}) P_{\tilde{U}^\perp} \mathbf{e}_k \\ &= \mathbf{e}_j^T X^+ \mathbf{E}[\mathbf{e} \mathbf{e}^T P_{U^\perp}] (X^+)^T (M^T P_{\tilde{U}^\perp}) \mathbf{e}_k + 0 \\ &\quad + l^2 \cdot \mathbf{e}_j^T (M^+ P_{\tilde{U}^\perp}) \mathbf{e}_k \\ &= 0 + 0 + 0 = 0 \end{aligned}$$

And as $\tilde{\beta}$ and $\tilde{\zeta}$ are Gaussians, having their covariance = 0 implies independence. \square

Having established that $\tilde{\beta}$ and $\tilde{\zeta}$ are independent Gaussians and specified their distributions, we continue with the proof of Theorem 3.1. We assume for now that there exists some small $a > 0$ s.t.

$$l^2(M^T M)_{j,j}^{-1} \leq \sigma^2(X^T X)_{j,j}^{-1} + l^2(M^T M)_{j,j}^{-1} \leq e^{2a} \cdot l^2(M^T M)_{j,j}^{-1} \quad (7)$$

Then, due to Corollary A.3, denoting the distributions $\mathcal{N}_1 = \mathcal{N}(0, l^2(M^T M)_{j,j}^{-1})$ and $\mathcal{N}_2 = \mathcal{N}(0, \sigma^2(X^T X)_{j,j}^{-1} + l^2(M^T M)_{j,j}^{-1})$, we have that for any $S \subset \mathbb{R}$ it holds that¹³

$$e^{-a} \Pr_{\tilde{\beta}_j \sim \mathcal{N}_1}[S] \leq \Pr_{\tilde{\beta}_j \sim \mathcal{N}_2}[S] \leq e^a \Pr_{\tilde{\beta}_j \sim \mathcal{N}_1}[S/e^a] \quad (8)$$

More specifically, denote the function

$$\begin{aligned} \tilde{t}(\psi, \|\boldsymbol{\xi}\|, \beta_j) &= \frac{\psi - \beta_j}{\|\boldsymbol{\xi}\| \sqrt{\frac{r}{r-p} (M^T M)_{j,j}^{-1}}} \\ &= \frac{\psi - \beta_j}{l \sqrt{(M^T M)_{j,j}^{-1}}} \bigg/ \frac{\|\boldsymbol{\xi}\| \sqrt{\frac{r}{r-p}}}{l} \end{aligned}$$

and observe that when we sample $\psi, \boldsymbol{\xi}$ independently s.t. $\psi \sim \mathcal{N}(\beta_j, l^2(M^T M)_{j,j}^{-1})$ and $\|\boldsymbol{\xi}\|^2 \sim \frac{l^2}{r} \chi_{r-p}^2$ then $\tilde{t}(\psi, \|\boldsymbol{\xi}\|, \beta_j)$ is distributed like a T -distribution with $r - p$

¹³In fact, it is possible to use standard techniques from differential privacy, and argue a similar result — that the probabilities of any event that depends on some function $f(\beta_j)$ under $\beta_j \sim \mathcal{N}_1$ and under $\beta_j \sim \mathcal{N}_2$ are close in the differential privacy sense.

degrees of freedom. And so, for any $\tau > 0$ we have that under such way to sample ψ, ξ we have $\Pr[\tilde{t}(\psi, \|\xi\|, \beta_j) > \tau] = 1 - \text{CDF}_{T_{r-p}}(\tau)$.

For any $\tau \geq 0$ and for any non-negative real value z let S_z^τ denote the suitable set of values s.t.

$$\begin{aligned} & \Pr \left\{ \begin{array}{l} \psi \sim \mathcal{N}(\beta_j, l^2(M^\top M)_{j,j}^{-1}) \\ \|\xi\|^2 \sim \frac{l^2}{r} \chi_{r-p}^2 \end{array} \right\} [\tilde{t}(\psi, \|\xi\|, \beta_j) > \tau] \\ &= \int_0^\infty \text{PDF}_{\frac{l^2}{r} \chi_{r-p}^2}(z) \cdot \Pr_{\{\psi \sim \mathcal{N}(0, l^2(M^\top M)_{j,j}^{-1})\}} [S_z^\tau] dz \end{aligned}$$

$$\text{That is, } S_z^\tau = \left(\tau \cdot z \sqrt{\frac{r}{r-p}(M^\top M)_{j,j}^{-1}}, \infty \right).$$

We now use Equation (8) (Since $\mathcal{N}(0, l^2(M^\top M)_{j,j}^{-1})$ is precisely \mathcal{N}_1) to deduce that

$$\begin{aligned} & \Pr \left\{ \begin{array}{l} \psi \sim \mathcal{N}(\beta_j, l^2(M^\top M)_{j,j}^{-1} + \sigma^2(X^\top X)_{j,j}^{-1}) \\ \|\xi\|^2 \sim \frac{l^2}{r} \chi_{r-p}^2 \end{array} \right\} [\tilde{t}(\psi, \|\xi\|, \beta_j) > \tau] \\ &= \int_0^\infty \text{PDF}_{\frac{l^2}{r} \chi_{r-p}^2}(z) \Pr_{\{\psi \sim \mathcal{N}(0, l^2(M^\top M)_{j,j}^{-1} + \sigma^2(X^\top X)_{j,j}^{-1})\}} [S_z^\tau] dz \\ &\leq e^a \int_0^\infty \text{PDF}_{\frac{l^2}{r} \chi_{r-p}^2}(z) \Pr_{\{\psi \sim \mathcal{N}(0, l^2(M^\top M)_{j,j}^{-1})\}} [S_z^\tau / e^a] dz \\ &\stackrel{(*)}{=} e^a \int_0^\infty \text{PDF}_{\frac{l^2}{r} \chi_{r-p}^2}(z) \Pr_{\{\psi \sim \mathcal{N}(0, l^2(M^\top M)_{j,j}^{-1})\}} [S_z^{\tau/e^a}] dz \\ &= e^a \Pr \left\{ \begin{array}{l} \psi \sim \mathcal{N}(\beta_j, l^2(M^\top M)_{j,j}^{-1}) \\ \|\xi\|^2 \sim \frac{l^2}{r} \chi_{r-p}^2 \end{array} \right\} [\tilde{t}(\psi, \|\xi\|, \beta_j) > \tau/e^a] \\ &= e^a (1 - \text{CDF}_{T_{r-p}}(\tau/e^a)) \end{aligned}$$

where the equality (*) follows from the fact that $S_z^\tau/c = S_{z/c}^\tau$ for any $c > 0$, since it is a non-negative interval. Analogously, we can also show that

$$\begin{aligned} & \Pr \left\{ \begin{array}{l} \psi \sim \mathcal{N}(\beta_j, l^2(M^\top M)_{j,j}^{-1} + \sigma^2(X^\top X)_{j,j}^{-1}) \\ \|\xi\|^2 \sim \frac{l^2}{r} \chi_{r-p}^2 \end{array} \right\} [\tilde{t}(\psi, \|\xi\|, \beta_j) > \tau] \\ &\geq e^{-a} \Pr \left\{ \begin{array}{l} \psi \sim \mathcal{N}(\beta_j, l^2(M^\top M)_{j,j}^{-1}) \\ \|\xi\|^2 \sim \frac{l^2}{r} \chi_{r-p}^2 \end{array} \right\} [\tilde{t}(\psi, \|\xi\|, \beta_j) > \tau] \\ &= e^{-a} (1 - \text{CDF}_{T_{r-p}}(\tau)) \end{aligned}$$

In other words, we have just shown that for any interval $I = (\tau, \infty)$ with $\tau \geq 0$ we have that

$$\Pr \left\{ \begin{array}{l} \psi \sim \mathcal{N}(\beta_j, l^2(M^\top M)_{j,j}^{-1} + \sigma^2(X^\top X)_{j,j}^{-1}) \\ \|\xi\|^2 \sim \frac{l^2}{r} \chi_{r-p}^2 \end{array} \right\} [\tilde{t}(\psi, \|\xi\|, \beta_j) \in I]$$

is lower bounded by $e^a \int_I \text{PDF}_{T_{r-p}}(z) dz$ and upper bounded by $e^a \int_{I/e^a} \text{PDF}_{T_{r-p}}(z) dz$. We can now repeat

the same argument for $I = (\tau_1, \tau_2)$ with $0 \leq \tau_1 < \tau_2$ (using an analogous definition of $S_z^{\tau_1, \tau_2}$), and again

for any $I = (\tau_1, \tau_2)$ with $\tau_1 < \tau_2 \leq 0$, and deduce that the PDF of the function $\tilde{t}(\psi, \|\xi\|, \beta_j)$ at x — where we sample $\psi \sim \mathcal{N}(\beta_j, l^2(M^\top M)_{j,j}^{-1} + \sigma^2(X^\top X)_{j,j}^{-1})$ and $\|\xi\|^2 \sim \frac{l^2}{r} \chi_{r-p}^2$ independently — lies in the range $(e^{-a} \text{PDF}_{T_{r-p}}(x), e^a \text{PDF}_{T_{r-p}}(x/e^a))$. And so, using Corollary B.4 and Claim B.5, we have that when $e \sim \mathcal{D}_l$, the distributions of $\tilde{\beta}_j$ and $\|\tilde{\zeta}\|^2$ are precisely as stated above, and so we have that the distribution of $\tilde{t}(\beta_j) \stackrel{\text{def}}{=} \tilde{t}(\tilde{\beta}_j, \|\tilde{\zeta}\|, \beta_j)$ has a PDF that at the point x is “sandwiched” between $e^{-a} \text{PDF}_{T_{r-p}}(x)$ and $e^a \text{PDF}_{T_{r-p}}(x/e^a)$.

Next, we aim to argue that this characterization of the PDF of $\tilde{t}(\beta_j)$ still holds when $e \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 I_{n \times n})$. It would be convenient to think of e as a sample in $\mathcal{N}(\mathbf{0}_n, \sigma^2 P_U) \times \mathcal{N}(\mathbf{0}_n, \sigma^2 P_{U^\perp})$. (So while in \mathcal{D}_l we have $P_U e \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 P_U)$ but $P_{U^\perp} e$ is fixed, now both $P_U e$ and $P_{U^\perp} e$ are sampled from spherical Gaussians.) The reason why the above still holds lies in the fact that $\tilde{t}(\beta_j)$ does not depend on l . In more details:

$$\begin{aligned} & \Pr_{e \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 I_{n \times n})} [\tilde{t}(\beta_j) \in I] \\ &= \int_{\mathbf{v}} \Pr_{e \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 I_{n \times n})} [\tilde{t}(\beta_j) \in I \mid P_{U^\perp} e = \mathbf{v}] \text{PDF}_{P_{U^\perp} e}(\mathbf{v}) d\mathbf{v} \\ &= \int_{\mathbf{v}} \Pr_{e \sim \mathcal{D}_l} [\tilde{t}(\beta_j) \in I \mid l = \|\mathbf{v}\|] \text{PDF}_{P_{U^\perp} e}(\mathbf{v}) d\mathbf{v} \\ &\leq \int_{\mathbf{v}} \left(e^a \int_{I/e^a} \text{PDF}_{T_{r-p}}(z) dz \right) \text{PDF}_{P_{U^\perp} e}(\mathbf{v}) d\mathbf{v} \\ &= \left(e^a \int_{I/e^a} \text{PDF}_{T_{r-p}}(z) dz \right) \int_{\mathbf{v}} \text{PDF}_{P_{U^\perp} e}(\mathbf{v}) d\mathbf{v} \\ &= e^a \int_{I/e^a} \text{PDF}_{T_{r-p}}(z) dz \end{aligned}$$

where the last transition is possible precisely because \tilde{t} is independent of l (or $\|\mathbf{v}\|$) — which is precisely what makes this t -value a pivot quantity. The proof of the lower bound is symmetric.

To conclude, we have shown that if Equation (7) holds, then for every interval $I \subset \mathbb{R}$ we have that $\Pr_{e \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 I_{n \times n})} [\tilde{t}(\beta_j) \in I]$ is lower bounded by $e^{-a} \Pr_{z \sim T_{r-p}} [z \in I]$ and upper bounded by $e^a \Pr_{z \sim T_{r-p}} [z \in (I/e^a)]$. So to conclude the proof of Theorem 3.1, we need to show that w.h.p such a as in Equation (7) exists.

Claim B.6. *In the homoscedastic model with Gaussian noise, if both n and r satisfy $n, r \geq p + \Omega(\log(1/\nu))$, then we have that $\sigma^2(X^\top X)_{j,j}^{-1} + l^2(M^\top M)_{j,j}^{-1} \geq l^2(M^\top M)_{j,j}^{-1}$ and*

$$\sigma^2(X^\top X)_{j,j}^{-1} + l^2(M^\top M)_{j,j}^{-1} \leq (1 + \frac{2(r-p)}{n-p}) \cdot l^2(M^\top M)_{j,j}^{-1}$$

Using $(1 + \frac{2(r-p)}{n-p}) \leq e^{\frac{2(r-p)}{n-p}}$, Theorem 3.1 now follows from plugging $a = \frac{r-p}{n-p}$ to our above discussion.

Proof. The lower bound is immediate from non-negativity of σ^2 and of $(X^\top X)_{j,j}^{-1} = \|(X^\top X)^{-1/2} \mathbf{e}_j\|^2$. We therefore prove the upper bound.

First, observe that $l^2 = \|P_{U^\perp} \mathbf{e}\|^2$ is sampled from $\sigma^2 \cdot \chi_{n-p}^2$ as U^\perp is of dimension $n-p$. Therefore, it holds that w.p. $\geq 1 - \nu/2$ that

$$\sigma^2 \left(\sqrt{n-p} - \sqrt{2 \ln(2/\nu)} \right)^2 \leq l^2$$

and assuming $n > p + 100 \ln(2/\nu)$ we therefore have $\sigma^2 \leq \frac{4}{3(n-p)} l^2$.

Secondly, we argue that when $r > p + 300 \ln(4/\nu)$ we have that w.p. $\geq 1 - \nu/2$ it holds that $\frac{3}{4} (X^\top X)_{j,j}^{-1} \leq (r-p) (X^\top R^\top R X)_{j,j}^{-1}$. To see this, first observe that by picking $R \sim \mathcal{N}(0_{r \times n}, I_{r \times r}, I_{n \times n})$ the distribution of the product $RX \sim \mathcal{N}(0_{r \times d}, I_{r \times r}, X^\top X)$ is identical to picking $Q \sim \mathcal{N}(0_{r \times d}, I_{r \times r}, I_{d \times d})$ and taking the product $Q(X^\top X)^{1/2}$. Therefore, the distribution of $(X^\top R^\top R X)^{-1}$ is identical to $((X^\top X)^{1/2} Q^\top Q (X^\top X)^{1/2})^{-1} = (X^\top X)^{-1/2} (Q^\top Q)^{-1} (X^\top X)^{-1/2}$. Denoting $\mathbf{v} = (X^\top X)^{-1/2} \mathbf{e}_j$ we have $\|\mathbf{v}\|^2 = (X^\top X)_{j,j}^{-1}$. Claim A.1 from (Sheffet, 2015) gives that w.p. $\geq 1 - \nu/2$ we have

$$\begin{aligned} (r-p) \cdot \mathbf{e}_j^\top \left((X^\top X)^{1/2} Q^\top Q (X^\top X)^{1/2} \right)^{-1} \mathbf{e}_j \\ = \mathbf{v}^\top \left(\frac{1}{r-p} Q^\top Q \right)^{-1} \mathbf{v} \geq \frac{3}{4} \mathbf{v}^\top \mathbf{v} = \frac{3}{4} (X^\top X)_{j,j}^{-1} \end{aligned}$$

which implies the required.

Combining the two inequalities we get:

$$\begin{aligned} \sigma^2 (X^\top X)_{j,j}^{-1} &\leq \frac{16l^2(r-p)}{n-p} (X^\top R^\top R X)_{j,j}^{-1} \\ &\leq \frac{2(r-p)}{n-p} l^2 (X^\top R^\top R X)_{j,j}^{-1} \end{aligned}$$

and as we denote $M = RX$ we are done. \square

We comment that our analysis in the proof of Claim B.6 implicitly assumes $r \ll n$ (as we do think of the projection R as dimensionality reduction), and so the ratio $\frac{r-p}{n-p}$ is small. However, a similar analysis holds for r which is comparable to n — in which we would argue that $\frac{\sigma^2 (X^\top X)_{j,j}^{-1} + l^2 (M^\top M)_{j,j}^{-1}}{\sigma^2 (X^\top X)^{-1}} \in [1, 1 + \eta]$ for some small η .

B.3. Proof of Theorem 3.3

Theorem B.7 (Theorem 3.3 restated.). *Fix a positive definite matrix $\Sigma \in \mathbb{R}^{p \times p}$. Fix parameters $\boldsymbol{\beta} \in \mathbb{R}^p$ and $\sigma^2 > 0$ and a coordinate j s.t. $\beta_j \neq 0$. Let X be a matrix whose n rows are sampled i.i.d from $\mathcal{N}(\mathbf{0}_p, \Sigma)$. Let \mathbf{y} be a vector s.t. $y_i - (X\boldsymbol{\beta})_i$ is sampled i.i.d from $\mathcal{N}(0, \sigma^2)$. Fix $\nu \in (0, 1/2)$ and $\alpha \in (0, 1/2)$. Then there exist constants C_1, C_2, C_3 and C_4 such that when we run Algorithm 1 over $[X; \mathbf{y}]$ with*

parameter r w.p. $\geq 1 - \nu$ we correctly α -reject the null hypothesis using \tilde{p}_0 (i.e., w.p. $\geq 1 - \nu$ Algorithm 1 returns matrix unaltered and we can estimate \tilde{t}_0 and verify that indeed $\tilde{p}_0 < \alpha \cdot e^{-\frac{r-p}{n-p}}$) provided

$$r \geq p + \max \left\{ C_1 \frac{\sigma^2 (\tilde{c}_\alpha^2 + \tilde{\tau}_\alpha^2)}{\beta_j^2 \sigma_{\min}(\Sigma)}, C_2 \ln(1/\nu) \right\}$$

and

$$n \geq \max \left\{ r, C_3 \frac{w^2}{\min\{\sigma_{\min}(\Sigma), \sigma^2\}}, C_4 (p + \ln(1/\nu)) \right\}$$

where $\tilde{c}_\alpha, \tilde{\tau}_\alpha$ denote the numbers s.t. $\int_{\tilde{c}_\alpha/e^{r-p}}^{\infty} \text{PDF}_{T_{r-p}}(x) dx = \frac{\alpha}{2} e^{-\frac{r-p}{n-p}}$ and $\int_{\tilde{\tau}_\alpha/e^{r-p}}^{\infty} \text{PDF}_{\mathcal{N}(0,1)}(x) dx = \frac{\alpha}{2} e^{-\frac{r-p}{n-p}}$ resp.

Proof. First we need to use the lower bound on n to show that indeed Algorithm 1 does not alter A , and that various quantities are not far from their expected values. Formally, we claim the following.

Proposition B.8. *Under the same lower bounds on n and r as in Theorem 3.3, w.p. $1 - \alpha - \nu$ we have that Theorem 3.1 holds and also that*

$$\|\tilde{\boldsymbol{\zeta}}\|^2 = \Theta\left(\frac{r-p}{r} \|P_{U^\perp} \mathbf{e}\|^2\right) = \Theta\left(\frac{r-p}{r} (n-p) \sigma^2\right)$$

and

$$(X^\top R^\top R X)_{j,j}^{-1} = \Theta\left(\frac{1}{r-p} (X^\top X)_{j,j}^{-1}\right)$$

Proof of Proposition B.8. First, we need to argue that we have enough samples as to have the gap $\sigma_{\min}^2([X; \mathbf{y}]) - w^2$ sufficiently large.

Since $\mathbf{x}_i \sim \mathcal{N}(0, \Sigma)$, and $y_i = \boldsymbol{\beta}^\top \mathbf{x}_i + e_i$ with $e_i \sim \mathcal{N}(0, \sigma^2)$, we have that the concatenation $(\mathbf{x}_i \circ y_i)$ is also sampled from a Gaussian. Clearly, $\mathbf{E}[y_i] = \boldsymbol{\beta}^\top \mathbf{E}[\mathbf{x}_i] + \mathbf{E}[e_i] = 0$. Similarly, $\mathbf{E}[x_{i,j} y_i] = \mathbf{E}[x_{i,j} \cdot (\boldsymbol{\beta}^\top \mathbf{x}_i + e_i)] = (\Sigma \boldsymbol{\beta})_j$ and $\mathbf{E}[y_i^2] = \mathbf{E}[e_i^2] + \mathbf{E}[\|X\boldsymbol{\beta}\|^2] = \sigma^2 + \mathbf{E}[\boldsymbol{\beta}^\top X^\top X \boldsymbol{\beta}] = \sigma^2 + \boldsymbol{\beta}^\top \Sigma \boldsymbol{\beta}$. Therefore, each row of A is an i.i.d sample of $\mathcal{N}(\mathbf{0}_{p+1}, \Sigma_A)$, with

$$\Sigma_A = \left(\begin{array}{c|c} \Sigma & \Sigma \boldsymbol{\beta} \\ \hline \boldsymbol{\beta}^\top \Sigma & \sigma^2 + \boldsymbol{\beta}^\top \Sigma \boldsymbol{\beta} \end{array} \right)$$

Denote $\lambda^2 = \sigma_{\min}(\Sigma)$. Then, to argue that $\sigma_{\min}(\Sigma_A)$ is large we use the lower bound from (Ma & Zarowski, 1995) (Theorem 3.1) combining with some simple arithmetic manipulations to deduce that $\sigma_{\min}(\Sigma_A) \geq \min\{\sigma_{\min}(\Sigma), \sigma^2\}$.

Having established a lower bound on $\sigma_{\min}(\Sigma_A)$, it follows that with $n = \Omega(p \ln(1/\nu))$ i.i.d draws from $\mathcal{N}(\mathbf{0}_{p+1}, \Sigma_A)$ we have w.p. $\leq \nu/4$ that $\sigma_{\min}(A^\top A) = o(n) \cdot \min\{\sigma_{\min}(\Sigma), \sigma^2\}$. Conditioned on $\sigma_{\min}(A^\top A) = \Omega(n\sigma_{\min}(\Sigma_A)) = \Omega(w^2)$ being large enough, we have that w.p. $\leq \nu/4$ over the randomness of Algorithm 1 the matrix A does not pass the if-condition and the output of the algorithm is not RA . Conditioned on Algorithm 1 outputting RA , and due to the lower bound $r = p + \Omega(\ln(1/\nu))$, we have that the result of Theorem 3.1 does not hold w.p. $\leq \alpha + \nu/4$. All in all we deduce that w.p. $\geq 1 - \alpha - 3\nu/4$ the result of Theorem 3.1 holds. And since we argue Theorem 3.1 holds, then the following two bounds that are used in the proof¹⁴ also hold:

$$(X^\top R^\top R X)_{j,j}^{-1} = \Theta\left(\frac{1}{r-p}(X^\top X)_{j,j}^{-1}\right)$$

$$\|P_{U^\perp} \mathbf{e}\|^2 = \Theta((n-p)\sigma^2)$$

Lastly, in the proof of Theorem 3.1 we argue that for a given $P_{U^\perp} \mathbf{e}$ the length $\|\tilde{\boldsymbol{\zeta}}\|^2$ is distributed like $\frac{\|P_{U^\perp} \mathbf{e}\|^2}{r} \chi_{r-p}^2$. Appealing again to the fact that $r = p + \Omega(\ln(1/\nu))$ we have that w.p. $\geq \nu/4$ it holds that $\|\tilde{\boldsymbol{\zeta}}\|^2 > 2(r-p) \frac{\|P_{U^\perp} \mathbf{e}\|^2}{r}$. Plugging in the value of $\|P_{U^\perp} \mathbf{e}\|^2$ concludes the proof of the proposition. \square

Based on Proposition B.8, we now show that we indeed reject the null-hypothesis (as we should). When Theorem 3.1 holds, reject the null-hypothesis iff $\tilde{p}_0 < \alpha \cdot e^{-\frac{r-p}{n-p}}$ which holds iff $|\tilde{t}_0| > e^{\frac{r-p}{n-p}} \tilde{\tau}_\alpha$. This implies we reject that null-hypothesis when $|\tilde{\beta}_j| > e^{\frac{r-p}{n-p}} \tilde{\tau}_\alpha \cdot \tilde{\sigma} \sqrt{(X^\top R^\top R X)_{j,j}^{-1}}$. Note that this bound is based on Corollary 3.2 that determines that $|\tilde{\beta}_j - \beta_j| = O\left(e^{\frac{r-p}{n-p}} \tilde{c}_\alpha \cdot \tilde{\sigma} \sqrt{(X^\top R^\top R X)_{j,j}^{-1}}\right)$. And so we have that w.p. $\geq 1 - \nu$ we α -reject the null hypothesis when it holds that $|\beta_j| > 3(\tilde{c}_\alpha + \tilde{\tau}_\alpha) \cdot \tilde{\sigma} \sqrt{(X^\top R^\top R X)_{j,j}^{-1}} \geq e^{\frac{r-p}{n-p}} (\tilde{c}_\alpha + \tilde{\tau}_\alpha) \tilde{\sigma} \sqrt{(X^\top R^\top R X)_{j,j}^{-1}}$ (due to the lower bound $n \geq r$).

Based on the bounds stated above we have that

$$\tilde{\sigma} = \|\tilde{\boldsymbol{\zeta}}\| \sqrt{\frac{r}{r-p}} = \Theta(\sigma \sqrt{n-p} \sqrt{\frac{r-p}{r}} \sqrt{\frac{r}{r-p}}) = \Theta(\sigma \sqrt{n-p})$$

and that

$$(X^\top R^\top R X)_{j,j}^{-1} = \Theta\left(\frac{1}{r-p}(X^\top X)_{j,j}^{-1}\right) = O\left(\frac{1}{r-p} \cdot \frac{1}{n\sigma_{\min}(\Sigma)}\right)$$

And so, a sufficient condition for rejecting the null-hypothesis is to have

$$|\beta_j| = \Omega\left((\tilde{c}_\alpha + \tilde{\tau}_\alpha) \sigma \sqrt{\frac{n-p}{r-p}} \cdot \sqrt{\frac{1}{n\sigma_{\min}(\Sigma)}}\right)$$

¹⁴More accurately, both are bounds shown in Claim B.6.

$$= \Omega\left(e^{\frac{r-p}{n-p}} (\tilde{c}_\alpha + \tilde{\tau}_\alpha) \tilde{\sigma} \sqrt{(X^\top R^\top R X)_{j,j}^{-1}}\right)$$

which, given the lower bound $r = p + \Omega\left(\frac{(\tilde{c}_\alpha + \tilde{\tau}_\alpha)^2 \sigma^2}{\beta_j^2 \sigma_{\min}(\Sigma)}\right)$ indeed holds. \square

C. Projected Ridge Regression

In this section we deal with the case that our matrix does not pass the if-condition of Algorithm 1. In this case, the matrix is appended with a $d \times d$ -matrix which is $wI_{d \times d}$.

Denoting $A' = \begin{bmatrix} A \\ w \cdot I_{d \times d} \end{bmatrix}$ we have that the algorithm's output is RA' .

Similarly to before, we are going to denote $d = p + 1$ and decompose $A = [X; \mathbf{y}]$ with $X \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^n$, with the standard assumption of $\mathbf{y} = X\boldsymbol{\beta} + \mathbf{e}$ and e_i sampled i.i.d from $\mathcal{N}(0, \sigma^2)$.¹⁵ We now need to introduce some additional notation. We denote the appended matrix and vectors X' and \mathbf{y}' s.t. $A' = [X'; \mathbf{y}']$. Meaning:

$$X' = \begin{bmatrix} X \\ wI_{p \times p} \\ \mathbf{0}_p^\top \end{bmatrix}$$

and

$$\mathbf{y}' = \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_p \\ w \end{bmatrix} = X'\boldsymbol{\beta} + \begin{bmatrix} \mathbf{e} \\ -w\boldsymbol{\beta} \\ w \end{bmatrix} \stackrel{\text{def}}{=} X'\boldsymbol{\beta} + \mathbf{e}'$$

And so we respectively denote $R = [R_1; R_2; R_3]$ with $R_1 \in \mathbb{R}^{r \times n}$, $R_2 \in \mathbb{R}^{r \times p}$ and $R_3 \in \mathbb{R}^{r \times 1}$ (so R_3 is a vector denoted as a matrix). Hence:

$$M' = RX' = R_1X + wR_2$$

and

$$R\mathbf{y}' = RX'\boldsymbol{\beta} + R\mathbf{e}' = R_1\mathbf{y} + wR_3 = R_1X\boldsymbol{\beta} + R_1\mathbf{e} + wR_3$$

And so, using the output RA' of Algorithm 1, we solve the linear regression problem derived from $\frac{1}{r}RX'$ and $\frac{1}{\sqrt{r}}R\mathbf{y}'$. I.e., we set

$$\boldsymbol{\beta}' = \arg \min_{\mathbf{z}} \frac{1}{r} \|R\mathbf{y}' - RX'\mathbf{z}\|^2$$

$$= (X'^\top R^\top R X')^{-1} (RX')^\top (R\mathbf{y}')$$

Sarlos' results (2006) regarding the Johnson Lindenstrauss transform give that, when R has sufficiently many rows, solving the latter optimization problem gives a good approximation for the solution of the optimization problem

$$\boldsymbol{\beta}^R = \arg \min_{\mathbf{z}} \|\mathbf{y}' - X'\mathbf{z}\|^2 = \arg \min_{\mathbf{z}} (\|\mathbf{y} - X\mathbf{z}\|^2 + w^2\|\mathbf{z}\|^2)$$

¹⁵Just as before, it is possible to denote any single column as \mathbf{y} and any subset of the remaining columns as X .

The latter problem is known as the Ridge Regression problem. Invented in the 60s (Tikhonov, 1963; Hoerl & Kennard, 1970), Ridge Regression is often motivated from the perspective of penalizing linear vectors whose coefficients are too large. It is also often applied in the case where X doesn't have full rank or is close to not having full-rank. That is because the Ridge Regression problem is always solvable. One can show that the minimizer $\beta^R = (X^T X + w^2 I_{p \times p})^{-1} X^T \mathbf{y}$ is the unique solution of the Ridge Regression problem and that the RHS is always defined (even when X is singular).

The original focus of Ridge Regression is on penalizing β^R for having large coefficients. Therefore, Ridge Regression actually poses a family of linear regression problems: $\min_{\mathbf{z}} \|y - X\mathbf{z}\| + \lambda \|\mathbf{z}\|^2$, where one may set λ to be any non-negative scalar. And so, much of the literature on Ridge Regression is devoted to the art of fine-tuning this penalty term — either empirically or based on the λ that yields the best risk: $\|\mathbf{E}[\beta^R] - \beta\|^2 + \text{Var}(\beta^R)$.¹⁶ Here we propose a fundamentally different approach for the choice of the normalization factor — we set it so that solution of the regression problem would satisfy (ϵ, δ) -differential privacy (by projecting the problem onto a lower dimension).

While the solution of the Ridge Regression problem might have smaller risk than the OLS solution, it is not known how to derive t -values and/or reject the null hypothesis under Ridge Regression (except for using X to manipulate β^R back into $\hat{\beta} = (X^T X)^{-1} X^T \mathbf{y}$ and relying on OLS). In fact, prior to our work there was no need for such analysis! For confidence intervals one could just use the standard OLS, because access to X and \mathbf{y} was given.

Therefore, much for the same reason, we are unable to derive t -values under projected Ridge Regression.¹⁷ Clearly, there are situations where such confidence bounds simply cannot be derived. (Consider for example the case where $X = 0_{n \times p}$ and \mathbf{y} is just i.i.d draws from $\mathcal{N}(0, \sigma^2)$, so obviously $[X; y]$ gives no information about β .) Nonetheless, under additional assumptions about the data, our work can give confidence intervals for β_j , and in the case where the interval doesn't intersect the origin — assure us that $\text{sign}(\beta'_j) = \text{sign}(\beta_j)$ w.h.p.

Clearly, Sarlos' work (2006) gives an upper bound on the distance $\|\beta' - \beta^R\|$. However, such distance bound doesn't come with the coordinate by coordinate confidence guarantee we would like to have. In fact, it is not even clear from Sarlos' work that $\mathbf{E}[\beta'] = \beta^R$ (though it is obvious to see that $\mathbf{E}[(X^T R^T R X)^{-1} \beta^R] = \mathbf{E}[(R X)^T R \mathbf{y}']$). Here,

¹⁶Ridge Regression, as opposed to OLS, does not yield an unbiased estimator. I.e., $\mathbf{E}[\beta^R] \neq \beta$.

¹⁷Note: The naïve approach of using RX' and $R\mathbf{y}'$ to interpolate RX and $R\mathbf{y}$ and then apply Theorem 3.1 using these estimations of RX and $R\mathbf{y}$ ignores the noise added from appending the matrix A into A' , and it is therefore bound to produce inaccurate estimations of the t -values.

we show that $\mathbf{E}[\beta'] = \hat{\beta}$ which, more often than not, does not equal β^R .

Comment about notation. Throughout this section we assume X is of full rank and so $(X^T X)^{-1}$ is well-defined. If X isn't full-rank, then one can simply replace any occurrence of $(X^T X)^{-1}$ with $X^+(X^+)^T$. This makes all our formulas well-defined in the general case.

C.1. Running OLS on the Projected Data

In this section, we analyze the projected Ridge Regression, under the assumption (for now) that \mathbf{e} is fixed. That is, for now we assume that the only source of randomness comes from picking the matrix $R = [R_1; R_2; R_3]$. As before, we analyze the distribution over β' (see Equation (9)), and the value of the function we optimize at β' . Denoting $M' = R X'$, we can formally express the estimators:

$$\beta' = (M'^T M')^{-1} M'^T R \mathbf{y}' \quad (9)$$

$$\zeta' = \frac{1}{\sqrt{r}} (R \mathbf{y}' - R X' \beta') \quad (10)$$

Claim C.1. Given that $\mathbf{y} = X\beta + \mathbf{e}$ for a fixed \mathbf{e} , and given X and $M' = R X' = R_1 X + w R_2$ we have that

$$\begin{aligned} \beta' &\sim \mathcal{N}\left(\beta + X^+ \mathbf{e}, \right. \\ &\quad \left. (w^2 (\|\beta + X^+ \mathbf{e}\|^2 + 1) + \|P_{U^\perp} \mathbf{e}\|^2) (M'^T M')^{-1}\right) \\ \zeta' &\sim \mathcal{N}\left(\mathbf{0}_r, \right. \\ &\quad \left. \frac{w^2 (\|\beta + X^+ \mathbf{e}\|^2 + 1) + \|P_{U^\perp} \mathbf{e}\|^2}{r} (I_{r \times r} - M' M'^+)\right) \end{aligned}$$

and furthermore, β' and ζ' are independent of one another.

Proof. First, we write β' and ζ' explicitly, based on \mathbf{e} and projection matrices:

$$\begin{aligned} \beta' &= (M'^T M')^{-1} M'^T R \mathbf{y}' \\ &= M'^+ (R_1 X) \beta + M'^+ (R_1 \mathbf{e} + w R_3) \\ \zeta' &= \frac{1}{\sqrt{r}} (R \mathbf{y}' - R X' \beta') \\ &= \frac{1}{\sqrt{r}} (I_{r \times r} - M' M'^+) R \mathbf{e}' \\ &= \frac{1}{\sqrt{r}} P_{U^\perp} (R_1 \mathbf{e} - w R_2 \beta + w R_3) \end{aligned}$$

with U' denoting $\text{colspan}(M')$ and $P_{U'^\perp}$ denoting the projection onto the subspace U'^\perp .

Again, we break \mathbf{e} into an orthogonal composition: $\mathbf{e} = P_U \mathbf{e} + P_{U^\perp} \mathbf{e}$ with $U = \text{colspan}(X)$ (hence $P_U = X X^+$) and $U^\perp = \text{colspan}(X)^\perp$. Therefore,

$$\begin{aligned} \beta' &= M'^+ (R_1 X) \beta + M'^+ (R_1 X X^+ \mathbf{e} + R_1 P_{U^\perp} \mathbf{e} + w R_3) \\ &= M'^+ (R_1 X) (\beta + X^+ \mathbf{e}) + M'^+ (R_1 P_{U^\perp} \mathbf{e} + w R_3) \end{aligned}$$

whereas ζ' is essentially

$$\begin{aligned} & \frac{1}{\sqrt{r}}(I_{r \times r} - M' M'^+)(R_1 X X^+ \mathbf{e} + R_1 P_{U^\perp} \mathbf{e} - w R_2 \boldsymbol{\beta} + w R_3) \\ & \stackrel{(*)}{=} \frac{1}{\sqrt{r}}(I_{r \times r} - M' M'^+) \cdot \\ & \quad (R_1 X X^+ \mathbf{e} + R_1 P_{U^\perp} \mathbf{e} + (M' - w R_2) \boldsymbol{\beta} + w R_3) \\ & = \frac{1}{\sqrt{r}}(I_{r \times r} - M' M'^+) \cdot \\ & \quad (R_1 X (\boldsymbol{\beta} + X^+ \mathbf{e}) + R_1 P_{U^\perp} \mathbf{e} + w R_3) \end{aligned}$$

where equality (*) holds because $(I - M' M'^+) M' \mathbf{v} = \mathbf{0}$ for any \mathbf{v} .

We now aim to describe the distribution of R given that we know X' and $M' = R X'$. Since

$$\begin{aligned} M' &= R_1 X + w R_2 + 0 \cdot R_3 = R_1 X (X^+ X) + w R_2 \\ &= (R_1 P_U) X + w R_2 \end{aligned}$$

then M' is independent of R_3 and independent of $R_1 P_{U^\perp}$. Therefore, given X and M' the induced distribution over R_3 remains $R_3 \sim \mathcal{N}(\mathbf{0}_r, I_{r \times r})$, and similarly, given X and M' we have $R_1 P_{U^\perp} \sim \mathcal{N}(0_{r \times n}, I_{r \times r}, P_{U^\perp})$ (rows remain independent from one another, and each row is distributed like a spherical Gaussian in $\text{colspan}(X)^\perp$). And so, we have that $R_1 X = R_1 P_U X = M' - w R_2$, which in turn implies:

$$R_1 X \sim \mathcal{N}(M', I_{r \times r}, w^2 \cdot I_{p \times p})$$

multiplying this random matrix with a vector, we get

$$R_1 X (\boldsymbol{\beta} + X^+ \mathbf{e}) \sim \mathcal{N}(M' \boldsymbol{\beta} + M' X^+ \mathbf{e}, w^2 \|\boldsymbol{\beta} + X^+ \mathbf{e}\|^2 I_{r \times r})$$

and multiplying this random vector with a matrix we get

$$M'^+ R_1 X (\boldsymbol{\beta} + X^+ \mathbf{e}) \sim \mathcal{N}(\boldsymbol{\beta} + X^+ \mathbf{e}, w^2 \|\boldsymbol{\beta} + X^+ \mathbf{e}\|^2 (M'^T M')^{-1})$$

I.e.,

$$M'^+ R_1 X (\boldsymbol{\beta} + X^+ \mathbf{e}) \sim \|\boldsymbol{\beta} + X^+ \mathbf{e}\| \cdot \mathcal{N}(\mathbf{u}, w^2 (M'^T M')^{-1})$$

where \mathbf{u} denotes a unit-length vector in the direction of $\boldsymbol{\beta} + X^+ \mathbf{e}$.

Similar to before we have

$$\begin{aligned} R P_{U^\perp} &\sim \mathcal{N}(0_{r \times n}, I_{r \times r}, P_{U^\perp}) \\ &\Rightarrow M'^+(R P_{U^\perp} \mathbf{e}) \sim \mathcal{N}(\mathbf{0}_d, \|P_{U^\perp} \mathbf{e}\|^2 (M'^T M')^{-1}) \\ w R_3 &\sim \mathcal{N}(\mathbf{0}_r, w^2 I_{r \times r}) \\ &\Rightarrow M'^+(w R_3) \sim \mathcal{N}(\mathbf{0}_d, w^2 (M'^+ M')^{-1}) \end{aligned}$$

Therefore, the distribution of $\boldsymbol{\beta}'$, which is the sum of the 3 independent Gaussians, is as required.

Also, $\zeta' = \frac{1}{\sqrt{r}} P_{U^\perp} (R_1 X (\boldsymbol{\beta} + X^+ \mathbf{e}) + R_1 P_{U^\perp} \mathbf{e} + w R_3)$ is the sum of 3 independent Gaussians, which implies its distribution is

$$\begin{aligned} & \mathcal{N}\left(\frac{1}{\sqrt{r}} P_{U^\perp} M' (\boldsymbol{\beta} + X^+ \mathbf{e}), \right. \\ & \quad \left. \frac{1}{r} (w^2 (\|\boldsymbol{\beta} + X^+ \mathbf{e}\|^2 + 1) + \|P_{U^\perp} \mathbf{e}\|^2) P_{U^\perp}\right) \end{aligned}$$

I.e., $\mathcal{N}(\mathbf{0}_r, \frac{1}{r} (w^2 (\|\boldsymbol{\beta} + X^+ \mathbf{e}\|^2 + 1) + \|P_{U^\perp} \mathbf{e}\|^2) P_{U^\perp})$ as $P_{U^\perp} M' = 0_{r \times r}$.

Finally, observe that $\boldsymbol{\beta}'$ and ζ' are independent as the former depends on the projection of the spherical Gaussian $R_1 X (\boldsymbol{\beta} + X^+ \mathbf{e}) + R_1 P_{U^\perp} \mathbf{e} + w R_3$ on U' , and the latter depends on the projection of the same multivariate Gaussian on U'^\perp . \square

Observe that Claim C.1 assumes \mathbf{e} is given. This may seem somewhat strange, since without assuming anything about \mathbf{e} there can be many combinations of $\boldsymbol{\beta}$ and \mathbf{e} for which $\mathbf{y} = X \boldsymbol{\beta} + \mathbf{e}$. However, we always have that $\boldsymbol{\beta} + X^+ \mathbf{e} = X^+ \mathbf{y} = \hat{\boldsymbol{\beta}}$. Similarly, it is always the case the $P_{U^\perp} \mathbf{e} = (I - X X^+) \mathbf{y} = \zeta$. (Recall OLS definitions of $\hat{\boldsymbol{\beta}}$ and ζ in Equation (1) and (2).) Therefore, the distribution of $\boldsymbol{\beta}'$ and ζ' is unique (once \mathbf{y} is set):

$$\begin{aligned} \boldsymbol{\beta}' &\sim \mathcal{N}\left(\hat{\boldsymbol{\beta}}, (w^2 (\|\hat{\boldsymbol{\beta}}\|^2 + 1) + \|\zeta\|^2) (M'^T M')^{-1}\right) \\ \zeta' &\sim \mathcal{N}\left(\mathbf{0}_r, \frac{w^2 (\|\hat{\boldsymbol{\beta}}\|^2 + 1) + \|\zeta\|^2}{r} (I_{r \times r} - M' M'^+)\right) \end{aligned}$$

And so for a given dataset $[X; \mathbf{y}]$ we have that $\boldsymbol{\beta}'$ serves as an approximation for $\hat{\boldsymbol{\beta}}$.

An immediate corollary of Claim C.1 is that for any fixed \mathbf{e} it holds that the quantity $t'(\beta_j) = \frac{\beta'_j - (\beta_j + (X^+ \mathbf{e})_j)}{\|\zeta'\| \sqrt{\frac{r}{r-p} \cdot (M'^T M')_{j,j}^{-1}}} = \frac{\beta'_j - \beta_j}{\|\zeta'\| \sqrt{\frac{r}{r-p} \cdot (M'^T M')_{j,j}^{-1}}}$ is distributed like a T_{r-p} -distribution. Therefore, the following theorem follows immediately.

Theorem C.2. Fix $X \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}$. Define $\hat{\boldsymbol{\beta}} = X^+ \mathbf{y}$ and $\zeta = (I - X X^+) \mathbf{y}$. Let $R X'$ and $R \mathbf{y}'$ denote the result of applying Algorithm 1 to the matrix $A = [X; \mathbf{y}]$ when the algorithm appends the data with a $w \cdot I$ matrix. Fix a coordinate j and any $\alpha \in (0, 1/2)$. When computing $\boldsymbol{\beta}'$ and ζ' as in Equations (9) it and (10), we have that w.p. $\geq 1 - \alpha$ it holds that

$$\hat{\beta}_j \in \left(\beta'_j \pm c'_\alpha \|\zeta'\| \sqrt{\frac{r}{r-p} \cdot (M'^T M')_{j,j}^{-1}} \right)$$

where c'_α denotes the number such that $(-c'_\alpha, c'_\alpha)$ contains $1 - \alpha$ mass of the T_{r-p} -distribution.

Note that Theorem C.2, much like the rest of the discussion in this Section, builds on \mathbf{y} being fixed, which means β'_j serves as an approximation for $\hat{\beta}_j$. Yet our goal is to argue about similarity (or proximity) between β'_j and β_j . To that end, we combine the standard OLS confidence interval — which says that w.p. $\geq 1 - \alpha$ over the randomness of picking \mathbf{e} in the homoscedastic model we have

$|\beta_j - \hat{\beta}_j| \leq c_\alpha \|\zeta\| \sqrt{\frac{(X^T X)_{j,j}^{-1}}{n-p}}$ — with the confidence interval of Theorem C.2 above, and deduce that w.p. $\geq 1 - \alpha$

we have that $|\beta'_j - \beta_j|$ is at most

$$O\left(c_\alpha \frac{\|\zeta\| \sqrt{(X^\top X)_{j,j}^{-1}}}{\sqrt{n-p}} + c'_\alpha \frac{\|\zeta'\| \sqrt{r(M'^\top M')_{j,j}^{-1}}}{\sqrt{r-p}}\right) \quad (11)$$

¹⁸And so, in the next section, our goal is to give conditions under which the interval of Equation (11) isn't much larger in comparison to the interval length of $c'_\alpha \frac{\|\zeta'\|}{\sqrt{r-p}} \sqrt{r(M'^\top M')_{j,j}^{-1}}$ we get from Theorem C.2; and more importantly — conditions that make the interval of Theorem C.2 useful and not too large. (Note, in expectation $\frac{\|\zeta'\|}{\sqrt{r-p}}$ is about $\sqrt{(w^2 + w^2\|\hat{\beta}\|^2 + \|\zeta\|^2)/r}$. So, for example, in situations where $\|\hat{\beta}\|$ is very large, this interval isn't likely to inform us as to the sign of β_j .)

Motivating Example. A good motivating example for the discussion in the following section is when $[X; \mathbf{y}]$ is a strict submatrix of the dataset A . That is, our data contains many variables for each entry (i.e., the dimensionality d of each entry is large), yet our regression is made only over a modest subset of variables out of the d . In this case, the least singular value of A might be too small, causing the algorithm to alter A ; however, $\sigma_{\min}(X^\top X)$ could be sufficiently large so that had we run Algorithm 1 only on $[X; \mathbf{y}]$ we would not alter the input. (Indeed, a differentially private way for finding a subset of the variables that induce a submatrix with high σ_{\min} is an interesting open question, partially answered — for a single regression — in the work of Thakurta and Smith (Thakurta & Smith, 2013).) Indeed, the conditions we specify in the following section depend on $\sigma_{\min}(\frac{1}{n}X^\top X)$, which, for a zero-mean data, the minimal variance of the data in any direction. For this motivating example, indeed such variance isn't necessarily small.

C.2. Conditions for Deriving a Confidence Interval for Ridge Regression

Looking at the interval specified in Equation (11), we now give an upper bound on the the random quantities in this interval: $\|\zeta\|$, $\|\zeta'\|$, and $(M'^\top M')_{j,j}^{-1}$. First, we give bound that are dependent on the randomness in R (i.e., we continue to view \mathbf{e} as fixed).

Proposition C.3. *For any $\nu \in (0, 1/2)$, if we have $r = p + \Omega(\ln(1/\nu))$ then with probability \geq*

¹⁸Observe that w.p. $\geq 1 - \alpha$ over the randomness of \mathbf{e} we have that $|\beta_j - \hat{\beta}_j| \leq c_\alpha \|\zeta\| \sqrt{\frac{(X^\top X)_{j,j}^{-1}}{n-p}}$, and w.p. $\geq 1 - \alpha$ over the randomness of R we have that $|\beta'_j - \hat{\beta}_j| \leq c'_\alpha \|\zeta'\| \sqrt{\frac{r}{r-p}} \cdot (M'^\top M')_{j,j}^{-1}$. So technically, to give a $(1 - \alpha)$ -confidence interval around β'_j that contains β_j w.p. $\geq 1 - \alpha$, we need to use $c_{\alpha/2}$ and $c'_{\alpha/2}$ instead of c_α and c'_α resp. To avoid overburdening the reader with what we already see as too many parameters, we switch to asymptotic notation.

$1 - \nu$ over the randomness of R we have $(r - p)(M'^\top M')_{j,j}^{-1} = \Theta((w^2 I_{p \times p} + X^\top X)_{j,j}^{-1})$ and $\frac{\|\zeta'\|^2}{r-p} = \Theta(\frac{w^2 + w^2\|\hat{\beta}\|^2 + \|\zeta\|^2}{r})$.

Proof. The former bound follows from known results on the Johnson-Lindenstrauss transform (as were shown in the proof of Claim B.6). The latter bound follows from standard concentration bounds of the χ^2 -distribution. \square

Plugging in the result of Proposition C.3 to Equation (11) we get that w.p. $\geq 1 - \nu$ the difference $|\beta'_j - \beta_j|$ is at most

$$O\left(c_\alpha \frac{\|\zeta\|}{\sqrt{n-p}} \sqrt{(X^\top X)_{j,j}^{-1}} + c'_\alpha \sqrt{\frac{w^2 + w^2\|\hat{\beta}\|^2 + \|\zeta\|^2}{r-p}} \sqrt{(w^2 I_{p \times p} + X^\top X)_{j,j}^{-1}}\right) \quad (12)$$

We will also use the following proposition.

Proposition C.4.

$$(X^\top X)_{j,j}^{-1} \leq \left(1 + \frac{w^2}{\sigma_{\min}(X^\top X)}\right) (w^2 I_{p \times p} + X^\top X)_{j,j}^{-1}$$

Proof. We have that

$$\begin{aligned} & (X^\top X)^{-1} \\ &= (X^\top X)^{-1} (X^\top X + w^2 I_{p \times p}) (X^\top X + w^2 I_{p \times p})^{-1} \\ &= (X^\top X + w^2 I_{p \times p})^{-1} + w^2 (X^\top X)^{-1} (X^\top X + w^2 I_{p \times p})^{-1} \\ &= (I_{p \times p} + w^2 (X^\top X)^{-1}) (X^\top X + w^2 I_{p \times p})^{-1} \\ &= (X^\top X + w^2 I_{p \times p})^{-1/2} \\ & \quad (I_{p \times p} + w^2 (X^\top X)^{-1}) \\ & \quad (X^\top X + w^2 I_{p \times p})^{-1/2} \end{aligned}$$

where the latter holds because $(I_{p \times p} + w^2 (X^\top X)^{-1})$ and $(X^\top X + w^2 I_{p \times p})^{-1}$ are diagonalizable by the same matrix V (the same matrix for which $(X^\top X) = VS^{-1}V^\top$). Since we have $\|I_{p \times p} + w^2 (X^\top X)^{-1}\| = 1 + \frac{w^2}{\sigma_{\min}^2(X)}$, it is clear that $(I_{p \times p} + w^2 (X^\top X)^{-1}) \preceq (1 + \frac{w^2}{\sigma_{\min}^2(X)}) I_{p \times p}$. We deduce that $(X^\top X)_{j,j}^{-1} = \mathbf{e}_j^\top (X^\top X)^{-1} \mathbf{e}_j \leq (1 + \frac{w^2}{\sigma_{\min}^2(X)}) (X^\top X + w^2 I_{p \times p})_{j,j}^{-1}$. \square

Based on Proposition C.4 we get from Equation (12) that

$|\beta'_j - \beta_j|$ is at most

$$O\left(c_\alpha \sqrt{\frac{\|\zeta\|^2(1 + \frac{w^2}{\sigma_{\min}(X^\top X)})}{n-p}} + c'_\alpha \sqrt{\frac{w^2 + w^2\|\hat{\beta}\|^2 + \|\zeta\|^2}{r-p}}\right) \sqrt{(w^2 I_{p \times p} + X^\top X)^{-1}_{j,j}} \quad (13)$$

And so, if it happens to be the case that exists some small $\eta > 0$ for which $\hat{\beta}, \zeta$ and w^2 satisfy

$$\frac{\|\zeta\|^2(1 + \frac{w^2}{\sigma_{\min}(X^\top X)})}{n-p} \leq \eta^2 \left(\frac{w^2 + w^2\|\hat{\beta}\|^2 + \|\zeta\|^2}{r-p} \right) \quad (14)$$

then we have that $\Pr[\beta_j \in (\beta'_j \pm O((1+\eta) \cdot c'_\alpha \|\zeta'\| \sqrt{\frac{r}{r-p} \cdot (M'^\top M')^{-1}_{j,j}}))] \geq 1 - \alpha$.¹⁹ Moreover, if in this case $|\beta_j| > c'_\alpha(1 + \eta) \sqrt{\frac{w^2 + w^2\|\hat{\beta}\|^2 + \|\zeta\|^2}{r-p}} \sqrt{(w^2 I_{p \times p} + X^\top X)^{-1}_{j,j}}$ then $\Pr[\text{sign}(\beta'_j) = \text{sign}(\beta_j)] \geq 1 - \alpha$. This is precisely what Claims C.5 and C.6 below do.

Claim C.5. *If there exists $\eta > 0$ s.t. $n - p \geq \frac{2}{\eta^2}(r - p)$ and $n^2 = \Omega\left(r^{3/2} \cdot \frac{B^2 \ln(1/\delta)}{\epsilon} \cdot \frac{1}{\eta^2 \sigma_{\min}(\frac{1}{n} X^\top X)}\right)$, then $\Pr[\beta_j \in (\beta'_j \pm O((1+\eta) \cdot c'_\alpha \|\zeta'\| \sqrt{\frac{r}{r-p} \cdot (M'^\top M')^{-1}_{j,j}}))] \geq 1 - \alpha$.*

Proof. Based on the above discussion, it is enough to argue that under the conditions of the claim, the constraint of Equation (14) holds. Since we require $\frac{\eta^2}{2} \geq \frac{r-p}{n-p}$ then it is evident that $\frac{\|\zeta\|^2}{n-p} \leq \frac{\eta^2 \|\zeta\|^2}{2(r-p)}$. So we now show that $\frac{\|\zeta\|^2}{n-p} \cdot \frac{w^2}{\sigma_{\min}(X^\top X)} \leq \frac{\eta^2 \|\zeta\|^2}{2(r-p)}$ under the conditions of the claim, and this will show the required. All that is left is some algebraic manipulations. It suffices to have:

$$\begin{aligned} \frac{\eta^2}{2} \cdot \frac{n-p}{r-p} \sigma_{\min}(X^\top X) &\geq \frac{\eta^2}{2} \cdot \frac{n^2}{r} \sigma_{\min}(\frac{1}{n} X^\top X) \\ &\geq \frac{32B^2 \sqrt{r} \ln(8/\delta)}{\epsilon} \geq w^2 \end{aligned}$$

which holds for $n^2 \geq r^{3/2} \cdot \frac{64B^2 \ln(1/\delta)}{\epsilon \eta^2} \sigma_{\min}(\frac{1}{n} X^\top X)^{-1}$, as we assume to hold. \square

Claim C.6. *Fix $\nu \in (0, \frac{1}{2})$. If (i) $n = p + \Omega(\ln(1/\nu))$, (ii) $\|\beta\|^2 = \Omega(\sigma^2 \|X^+\|_F^2 \ln(\frac{p}{\nu}))$ and (iii) $r - p = \Omega\left(\frac{(c'_\alpha)^2 (1+\eta)^2}{\beta_j^2} \left(1 + \|\beta\|^2 + \frac{\sigma^2}{\sigma_{\min}(\frac{1}{n} X^\top X)}\right)\right)$, then in the homoscedastic model, with probability $\geq 1 - \nu - \alpha$ we have that $\text{sign}(\beta_j) = \text{sign}(\beta'_j)$.*

¹⁹We assume $n \geq r$ so $c_\alpha < c'_\alpha$ as the T_{n-p} -distribution is closer to a normal Gaussian than the T_{r-p} -distribution.

Proof. Based on the above discussion, we aim to show that in the homoscedastic model (where each coordinate $e_i \sim \mathcal{N}(0, \sigma^2)$ independently) w.p. $\geq 1 - \nu$ it holds that the magnitude of β_j is greater than

$$c'_\alpha(1+\eta) \sqrt{\frac{w^2 + w^2\|\hat{\beta}\|^2 + \|\zeta\|^2}{r-p}} \sqrt{(w^2 I_{p \times p} + X^\top X)^{-1}_{j,j}}$$

To show this, we invoke Claim A.4 to argue that w.p. $\geq 1 - \nu$ we have (i) $\|\zeta\|^2 \leq 2\sigma^2(n-p)$ (since $n = p + \Omega(\ln(1/\nu))$), and (ii) $\|\hat{\beta}\|^2 \leq 2\|\beta\|^2$ (since $\|\beta - \hat{\beta}\|^2 \leq \sigma^2 \|X^+\|_F^2 \ln(\frac{p}{\nu})$ whereas $\|\beta\|^2 = \Omega(\sigma^2 \|X^+\|_F^2 \ln(\frac{p}{\nu}))$). We also use the fact that $(w^2 I_{p \times p} + X^\top X)^{-1}_{j,j} \leq (w^2 + \sigma_{\min}^{-1}(X^\top X))$, and then deduce that

$$\begin{aligned} (1+\eta)c'_\alpha \sqrt{\frac{w^2 + w^2\|\hat{\beta}\|^2 + \|\zeta\|^2}{r-p}} \sqrt{(w^2 I_{p \times p} + X^\top X)^{-1}_{j,j}} \\ \leq \frac{(1+\eta)c'_\alpha}{\sqrt{r-p}} \sqrt{2 \frac{w^2(1 + \|\beta\|^2) + \sigma^2(n-p)}{w^2 + \sigma_{\min}(X^\top X)}} \\ \leq \frac{(1+\eta)c'_\alpha}{\sqrt{r-p}} \sqrt{2(1 + \|\beta\|^2) + \frac{2\sigma^2(n-p)}{\sigma_{\min}(X^\top X)}} \leq |\beta_j| \end{aligned}$$

due to our requirement on $r - p$. \square

Observe, out of the 3 conditions specified in Claim C.6, condition (i) merely guarantees that the sample is large enough to argue that estimations are close to their expect value; and condition (ii) is there merely to guarantee that $\|\hat{\beta}\| \approx \|\beta\|$. It is condition (iii) which is non-trivial to hold, especially together with the conditions of Claim C.5 that pose other constraints in regards to r, n, η and the various other parameters in play. It is interesting to compare the requirements on r to the lower bound we get in Theorem 3.3 — especially the latter bound. The two bounds are strikingly similar, with the exception that here we also require $r - p$ to be greater than $\frac{1 + \|\beta\|^2}{\beta_j^2}$. This is part of the unfortunate effect of altering the matrix A : we cannot give confidence bounds only for the coordinates j for which β_j^2 is very small relative to $\|\beta\|^2$.

In summary, we require to have $n = p + \Omega(\ln(1/\nu))$ and that X contains enough sample points to have $\|\hat{\beta}\|$ comparable to $\|\beta\|$, and then set r and η such that (it is convenient to think of η as a small constant, say, $\eta = 0.1$)

- $r - p = O(\eta^2(n - p))$ (which implies $r = O(n)$)
- $r = O\left(\left(\eta^2 \frac{\epsilon n^2}{B^2 \ln(1/\delta)} \sigma_{\min}(\frac{1}{n} X^\top X)\right)^{\frac{2}{3}}\right)$
- $r - p = \Omega\left(\frac{1 + \|\beta\|^2}{\beta_j^2} + \frac{\sigma^2}{\beta_j^2} \cdot \sigma_{\min}^{-1}(\frac{1}{n} X^\top X)\right)$

to have that the $(1 - \alpha)$ -confidence interval around β_j^l does not intersect the origin. Once again, we comment that these conditions are sufficient but not necessary, and furthermore — even with these conditions holding — we do not make any claims of optimality of our confidence bound. That is because from Proposition C.4 onwards our discussion uses upper bounds that do not have corresponding lower bounds, to the best of our knowledge.

D. Confidence Intervals for “Analyze Gauss” Algorithm

To complete the picture, we now analyze the “Analyze Gauss” algorithm of Dwork et al (Dwork et al., 2014). Algorithm 2 works by adding random Gaussian noise to $A^T A$, where the noise is symmetric with each coordinate above the diagonal sampled i.i.d from $\mathcal{N}(0, \Delta^2)$ with $\Delta^2 = O\left(B^4 \frac{\log(1/\delta)}{\epsilon^2}\right)$.²⁰ Using the same notation for a sub-matrix of A as $[X; \mathbf{y}]$ as before, with $X \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^n$, we denote the output of Algorithm 2 as

$$\left(\begin{array}{c|c} \widetilde{X^T X} & \widetilde{X^T \mathbf{y}} \\ \hline \widetilde{\mathbf{y}^T X} & \widetilde{\mathbf{y}^T \mathbf{y}} \end{array} \right) = \left(\begin{array}{c|c} X^T X + N & X^T \mathbf{y} + \mathbf{n} \\ \hline \mathbf{y}^T X + \mathbf{n}^T & \mathbf{y}^T \mathbf{y} + m \end{array} \right) \quad (15)$$

where N is a symmetric $p \times p$ -matrix, \mathbf{n} is a p -dimensional vector and m is a scalar, whose coordinates are sampled i.i.d from $\mathcal{N}(0, \Delta^2)$.

Using the output of Algorithm 2, it is simple to derive analogues of $\hat{\beta}$ and $\|\zeta\|^2$ (Equations (1) and (2))

$$\tilde{\beta} = \left(\widetilde{X^T X} \right)^{-1} \widetilde{X^T \mathbf{y}} = (X^T X + N)^{-1} (X^T \mathbf{y} + \mathbf{n}) \quad (16)$$

$$\begin{aligned} \|\widetilde{\zeta}\|^2 &= \widetilde{\mathbf{y}^T \mathbf{y}} - 2 \widetilde{\mathbf{y}^T X} \tilde{\beta} + \tilde{\beta}^T \widetilde{X^T X} \tilde{\beta} \\ &= \widetilde{\mathbf{y}^T \mathbf{y}} - \widetilde{\mathbf{y}^T X} \widetilde{X^T X}^{-1} \widetilde{X^T \mathbf{y}} \end{aligned} \quad (17)$$

We now argue that it is possible to use $\tilde{\beta}_j$ and $\|\widetilde{\zeta}\|^2$ to get a confidence interval for β_j under certain conditions.

Theorem D.1. Fix $\alpha, \nu \in (0, \frac{1}{2})$. Assume that there exists $\eta \in (0, \frac{1}{2})$ s.t. $\sigma_{\min}(X^T X) > \Delta \sqrt{p \ln(1/\nu)}/\eta$. Under the homoscedastic model, given β and σ^2 , if we assume also that $\|\beta\| \leq B$ and $\|\hat{\beta}\| = \|(X^T X)^{-1} X^T \mathbf{y}\| \leq B$, then w.p. $\geq 1 - \alpha - \nu$ it holds that $|\beta_j - \tilde{\beta}_j|$ it at most

$$O\left(\rho \cdot \sqrt{\left(\widetilde{X^T X}_{j,j}^{-1} + \Delta \sqrt{p \ln(1/\nu)} \cdot \widetilde{X^T X}_{j,j}^{-2}\right) \ln(1/\alpha)}\right)$$

²⁰It is easy to see that the l_2 -global sensitivity of the mapping $A \mapsto A^T A$ is $\propto B^4$. Fix any A_1, A_2 that differ on one row which is some vector \mathbf{v} with $\|\mathbf{v}\| = B$ in A_1 and the all zero vector in A_2 . Then $GS_2^2 = \|A_1^T A_1 - A_2^T A_2\|_F^2 = \|\mathbf{v}\mathbf{v}^T\|_F^2 = \text{trace}(\mathbf{v}\mathbf{v}^T \cdot \mathbf{v}\mathbf{v}^T) = (\mathbf{v}^T \mathbf{v})^2 = B^4$.

$$+ \Delta \sqrt{\widetilde{X^T X}_{j,j}^{-2} \cdot \ln(1/\nu)} \cdot (B\sqrt{p} + 1)$$

where ρ is such that ρ^2 is w.h.p an upper bound on σ^2 , defined as

$$\rho^2 \stackrel{\text{def}}{=} \left(\frac{1}{\sqrt{n-p-2\sqrt{\ln(4/\alpha)}}} \right)^2 \cdot \left(\|\widetilde{\zeta}\|^2 - C \cdot \left(\Delta \frac{B^2 \sqrt{p}}{1-\eta} \sqrt{\ln(1/\nu)} + \Delta^2 \|\widetilde{X^T X}^{-1}\|_F \cdot \ln(p/\nu) \right) \right)$$

for some large constant C .

We comment that in practice, instead of using ρ , it might be better to use the MLE of σ^2 , namely:

$$\overline{\sigma^2} \stackrel{\text{def}}{=} \frac{1}{n-p} \left(\|\widetilde{\zeta}\|^2 + \Delta^2 \|\widetilde{X^T X}^{-1}\|_F \right)$$

instead of ρ^2 , the upper bound we derived for σ^2 . (Replacing an unknown variable with its MLE estimator is a common approach in applied statistics.) Note that the assumption that $\|\beta\| \leq B$ is fairly benign once we assume each row has bounded l_2 -norm. The assumption $\|\hat{\beta}\| \leq B$ simply assumes that $\hat{\beta}$ is a reasonable estimation of β , which is likely to hold if we assume that $X^T X$ is well-spread. The assumption about the magnitude of the least singular value of $X^T X$ is therefore the major one. Nonetheless, in the case we considered before where each row in X is sampled i.i.d from $\mathcal{N}(\mathbf{0}, \Sigma)$, this assumption merely means that n is large enough s.t. $n = \tilde{\Omega}\left(\frac{\Delta \sqrt{p \ln(1/\nu)}}{\eta \cdot \sigma_{\min}(\Sigma)}\right)$.

In order to prove Theorem D.1, we require the following proposition.

Proposition D.2. Fix any $\nu \in (0, \frac{1}{2})$. Fix any matrix $M \in \mathbb{R}^{p \times p}$. Let $\mathbf{v} \in \mathbb{R}^p$ be a vector with each coordinate sampled independently from a Gaussian $\mathcal{N}(0, \Delta^2)$. Then we have that $\Pr \left[\|M\mathbf{v}\| > \Delta \cdot \|M\|_F \sqrt{2 \ln(2p/\nu)} \right] < \nu$.

Proof. Given M , we have that $M\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \Delta^2 \cdot MM^T)$. Denoting M 's singular values as sv_1, \dots, sv_p , we can rotate $M\mathbf{v}$ without affecting its l_2 -norm and infer that $\|M\mathbf{v}\|^2$ is distributed like a sum on p independent Gaussians, each sampled from $\mathcal{N}(0, \Delta^2 \cdot sv_i^2)$. Standard union bound gives that w.p. $\geq 1 - \nu$ non of the p Gaussians exceeds its standard deviation by a factor of $\sqrt{2 \ln(2p/\nu)}$. Hence, w.p. $\geq 1 - \nu$ it holds that $\|M\mathbf{v}\|^2 \leq 2\Delta^2 \sum_i sv_i^2 \ln(2p/\nu) = 2\Delta^2 \cdot \text{trace}(MM^T) \cdot \ln(2p/\nu)$. \square

Our proof also requires the use of the following equality, that holds for any invertible A and any matrix B s.t. $I + B \cdot A^{-1}$ is invertible:

$$(A + B)^{-1} = A^{-1} - A^{-1} (I + BA^{-1})^{-1} BA^{-1}$$

In our case, we have

$$\begin{aligned}
 & \widetilde{X^\top X}^{-1} \\
 &= (X^\top X + N)^{-1} \\
 &= (X^\top X)^{-1} - (X^\top X)^{-1} (I + N(X^\top X)^{-1})^{-1} N (X^\top X)^{-1} \\
 &= (X^\top X)^{-1} \left(I - (I + N(X^\top X)^{-1})^{-1} N (X^\top X)^{-1} \right) \\
 &\stackrel{\text{def}}{=} (X^\top X)^{-1} (I - Z \cdot (X^\top X)^{-1}) \tag{18}
 \end{aligned}$$

Proof of Theorem D.1. Fix $\nu > 0$. First, we apply to standard results about Gaussian matrices, such as (Tao, 2012) (used also by (Dwork et al., 2014) in their analysis), to see that w.p. $\geq 1 - \nu/6$ we have $\|N\| = O(\Delta\sqrt{p\ln(1/\nu)})$. And so, for the remainder of the proof we fix N subject to having bounded operator norm. Note that by fixing N we fix $\widetilde{X^\top X}$.

Recall that in the homoscedastic model, $\mathbf{y} = X\boldsymbol{\beta} + \mathbf{e}$ with each coordinate of \mathbf{e} sampled i.i.d from $\mathcal{N}(0, \sigma^2)$. We therefore have that

$$\begin{aligned}
 \widetilde{\boldsymbol{\beta}} &= \widetilde{X^\top X}^{-1} (X^\top \mathbf{y} + \mathbf{n}) = \widetilde{X^\top X}^{-1} (X^\top X \boldsymbol{\beta} + X^\top \mathbf{e} + \mathbf{n}) \\
 &= \widetilde{X^\top X}^{-1} (\widetilde{X^\top X} - N) \boldsymbol{\beta} + \widetilde{X^\top X}^{-1} X^\top \mathbf{e} + \widetilde{X^\top X}^{-1} \mathbf{n} \\
 &= \boldsymbol{\beta} - \widetilde{X^\top X}^{-1} N \boldsymbol{\beta} + \widetilde{X^\top X}^{-1} X^\top \mathbf{e} + \widetilde{X^\top X}^{-1} \mathbf{n}
 \end{aligned}$$

Denoting the j -th row of $\widetilde{X^\top X}^{-1}$ as $\widetilde{X^\top X}_{j \rightarrow}^{-1}$ we deduce:

$$\widetilde{\beta}_j = \beta_j - \widetilde{X^\top X}_{j \rightarrow}^{-1} N \boldsymbol{\beta} + \widetilde{X^\top X}_{j \rightarrow}^{-1} X^\top \mathbf{e} + \widetilde{X^\top X}_{j \rightarrow}^{-1} \mathbf{n} \tag{19}$$

We naïvely bound the size of the term $\widetilde{X^\top X}_{j \rightarrow}^{-1} N \boldsymbol{\beta}$ by $\left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\| \|N\| \|\boldsymbol{\beta}\| = O\left(\left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\| \cdot B \Delta \sqrt{p \ln(1/\nu)}\right)$.

To bound $\widetilde{X^\top X}_{j \rightarrow}^{-1} X^\top \mathbf{e}$ note that \mathbf{e} is chosen independently of $\widetilde{X^\top X}$ and since $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$ we have $\widetilde{X^\top X}_{j \rightarrow}^{-1} X^\top \mathbf{e} \sim \mathcal{N}\left(\mathbf{0}, \sigma^2 \cdot \mathbf{e}_j^\top \widetilde{X^\top X}^{-1} \cdot X^\top X \cdot \widetilde{X^\top X}^{-1} \mathbf{e}_j\right)$. Since we have

$$\begin{aligned}
 & \widetilde{X^\top X}^{-1} \cdot X^\top X \cdot \widetilde{X^\top X}^{-1} \\
 &= \widetilde{X^\top X}^{-1} \cdot (\widetilde{X^\top X} - N) \cdot \widetilde{X^\top X}^{-1} \\
 &= \widetilde{X^\top X}^{-1} - \widetilde{X^\top X}^{-1} \cdot N \cdot \widetilde{X^\top X}^{-1}
 \end{aligned}$$

we can bound the variance of $\widetilde{X^\top X}_{j \rightarrow}^{-1} X^\top \mathbf{e}$ by $\sigma^2 \left(\widetilde{X^\top X}_{j,j}^{-1} + \|N\| \cdot \left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\|^2 \right)$. Appealing to Gaussian concentration bounds, we have that w.p. $\geq 1 - \alpha/2$ the absolute value of this Gaussian is at most $O\left(\sqrt{\left(\widetilde{X^\top X}_{j,j}^{-1} + \Delta\sqrt{p\ln(1/\nu)} \cdot \left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\|^2\right) \sigma^2 \ln(1/\alpha)}\right)$.

To bound $\widetilde{X^\top X}_{j \rightarrow}^{-1} \mathbf{n}$ note that $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \Delta^2 I)$ is sampled independently of $\widetilde{X^\top X}$. We therefore have that $\widetilde{X^\top X}_{j \rightarrow}^{-1} \mathbf{n} \sim \mathcal{N}(0, \Delta^2 \left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\|^2)$. Gaussian concentration bounds give that w.p. $\geq 1 - \nu/6$ we have $|\widetilde{X^\top X}_{j \rightarrow}^{-1} \mathbf{n}| = O\left(\Delta \left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\| \sqrt{\ln(1/\nu)}\right)$.

Plugging this into our above bounds on all terms that appear in Equation (19) we have that w.p. $\geq 1 - \nu/2 - \alpha/2$ we have that $|\widetilde{\beta}_j - \beta_j|$ is at most

$$\begin{aligned}
 & O\left(\left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\| \cdot B \Delta \sqrt{p \ln(1/\nu)}\right) \\
 & + O\left(\sigma \sqrt{\left(\widetilde{X^\top X}_{j,j}^{-1} + \Delta\sqrt{p\ln(1/\nu)} \cdot \left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\|^2\right) \ln(1/\alpha)}\right) \\
 & + O\left(\Delta \left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\| \sqrt{\ln(1/\nu)}\right)
 \end{aligned}$$

Note that due to the symmetry of $\widetilde{X^\top X}$ we have $\left\| \widetilde{X^\top X}_{j \rightarrow}^{-1} \right\|^2 = \widetilde{X^\top X}_{j,j}^{-2}$ (the (j, j) -coordinate of the matrix $\widetilde{X^\top X}^{-2}$), thus $|\widetilde{\beta}_j - \beta_j|$ is at most

$$\begin{aligned}
 & O\left(\sigma \cdot \sqrt{\left(\widetilde{X^\top X}_{j,j}^{-1} + \Delta\sqrt{p\ln(1/\nu)} \cdot \widetilde{X^\top X}_{j,j}^{-2}\right) \ln(1/\alpha)}\right) \\
 & + \Delta \sqrt{\widetilde{X^\top X}_{j,j}^{-2} \cdot \ln(1/\nu) \cdot (B\sqrt{p} + 1)} \tag{20}
 \end{aligned}$$

All of the terms appearing in Equation (20) are known given $\widetilde{X^\top X}$, except for σ — which is a parameter of the model. Next, we derive an upper bound on σ which we can then plug into Equation (20) to complete the proof of the theorem and derive a confidence interval for β_j .

Recall Equation (17), according to which we have

$$\begin{aligned}
 \|\widehat{\zeta}\|^2 &= \widehat{\mathbf{y}}^\top \widehat{\mathbf{y}} - \widehat{\mathbf{y}}^\top \widehat{X} \widehat{X}^\top \widehat{X}^{-1} \widehat{X}^\top \widehat{\mathbf{y}} \\
 &\stackrel{(18)}{=} \mathbf{y}^\top \mathbf{y} + m \\
 &\quad - (\mathbf{y}^\top X + \mathbf{n}^\top)(X^\top X)^{-1}(I - Z \cdot (X^\top X)^{-1})(X^\top \mathbf{y} + \mathbf{n}) \\
 &= \mathbf{y}^\top \mathbf{y} + m \\
 &\quad - \mathbf{y}^\top X (X^\top X)^{-1} X^\top \mathbf{y} \\
 &\quad + \mathbf{y}^\top X (X^\top X)^{-1} Z (X^\top X)^{-1} X^\top \mathbf{y} \\
 &\quad - 2\mathbf{y}^\top X (X^\top X)^{-1} \mathbf{n} \\
 &\quad + 2\mathbf{y}^\top X (X^\top X)^{-1} Z (X^\top X)^{-1} \mathbf{n} \\
 &\quad - \mathbf{n}^\top (X^\top X)^{-1} (I - Z \cdot (X^\top X)^{-1}) \mathbf{n}
 \end{aligned}$$

Recall that $\widehat{\beta} = (X^\top X)^{-1} X^\top \mathbf{y}$, and so we have

$$\begin{aligned}
 &= \mathbf{y}^\top (I - X(X^\top X)^{-1} X^\top) \mathbf{y} + m - \widehat{\beta}^\top Z \widehat{\beta} \\
 &\quad - 2\widehat{\beta}^\top (I - Z(X^\top X)^{-1}) \mathbf{n} - \mathbf{n}^\top \widehat{X}^\top \widehat{X}^{-1} \mathbf{n} \quad (21)
 \end{aligned}$$

and of course, both \mathbf{n} and m are chosen independently of $\widehat{X}^\top \widehat{X}$ and \mathbf{y} .

Before we bound each term in Equation (21), we first give a bound on $\|Z\|$. Recall, $Z = (I + N(X^\top X)^{-1})^{-1} N$. Recall our assumption (given in the statement of Theorem D.1) that $\sigma_{\min}(X^\top X) \geq \frac{\Delta}{\eta} \sqrt{p \ln(1/\nu)}$. This implies that $\|N(X^\top X)^{-1}\| \leq \|N\| \cdot \sigma_{\min}(X^\top X)^{-1} = O(\eta)$. Hence

$$\|Z\| \leq (\|I + N(X^\top X)^{-1}\|)^{-1} \cdot \|N\| = O\left(\frac{\Delta \sqrt{p \ln(1/\nu)}}{1-\eta}\right)$$

Moreover, this implies that $\|Z(X^\top X)^{-1}\| \leq O\left(\frac{\eta}{1-\eta}\right)$

and that $\|I - Z(X^\top X)^{-1}\| \leq O\left(\frac{1}{1-\eta}\right)$.

Armed with these bounds on the operator norms of Z and $(I - Z(X^\top X)^{-1})$ we bound the magnitude of the different terms in Equation (21).

- The term $\mathbf{y}^\top (I - X X^\top) \mathbf{y}$ is the exact term from the standard OLS, and we know it is distributed like $\sigma^2 \cdot \chi_{n-p}^2$ distribution. Therefore, it is greater than $\sigma^2(\sqrt{n-p} - 2\sqrt{\ln(4/\alpha)})^2$ w.p. $\geq 1 - \alpha/2$.
- The scalar m sampled from $m \sim \mathcal{N}(0, \Delta^2)$ is bounded by $O(\Delta \sqrt{\ln(1/\nu)})$ w.p. $\geq 1 - \nu/8$.
- Since we assume $\|\widehat{\beta}\| \leq B$, the term $\widehat{\beta}^\top Z \widehat{\beta}$ is upper bounded by $B^2 \|Z\| = O\left(\frac{B^2 \Delta \sqrt{p \ln(1/\nu)}}{1-\eta}\right)$.
- Denote $\mathbf{z}^\top \mathbf{n} = 2\widehat{\beta}^\top (I - Z(X^\top X)^{-1}) \mathbf{n}$. We thus have that $\mathbf{z}^\top \mathbf{n} \sim \mathcal{N}(0, \Delta^2 \|\mathbf{z}\|^2)$ and that its magnitude is at

most $O(\Delta \cdot \|\mathbf{z}\| \sqrt{\ln(1/\nu)})$ w.p. $\geq 1 - \nu/8$. We can upper bound $\|\mathbf{z}\| \leq 2\|\widehat{\beta}\| \|I - Z(X^\top X)^{-1}\| = O\left(\frac{B}{1-\eta}\right)$, and so this term's magnitude is upper bounded by $O\left(\frac{\Delta \cdot B \sqrt{\ln(1/\nu)}}{1-\eta}\right)$.

- Given our assumption about the least singular value of $X^\top X$ and with the bound on $\|N\|$, we have that $\sigma_{\min}(\widehat{X}^\top \widehat{X}) \geq \sigma_{\min}(X^\top X) - \|N\| > 0$ and so the symmetric matrix $\widehat{X}^\top \widehat{X}$ is a PSD. Therefore, the term $\mathbf{n}^\top \widehat{X}^\top \widehat{X}^{-1} \mathbf{n} = \|\widehat{X}^\top \widehat{X}^{-1/2} \mathbf{n}\|^2$ is strictly positive. Applying Proposition D.2 we have that w.p. $\geq 1 - \nu/8$ it holds that $\mathbf{n}^\top \widehat{X}^\top \widehat{X}^{-1} \mathbf{n} \leq O\left(\Delta^2 \|\widehat{X}^\top \widehat{X}^{-1}\|_F \cdot \ln(p/\nu)\right)$.

Plugging all of the above bounds into Equation (21) we get that w.p. $\geq 1 - \nu/2 - \alpha/2$ it holds that

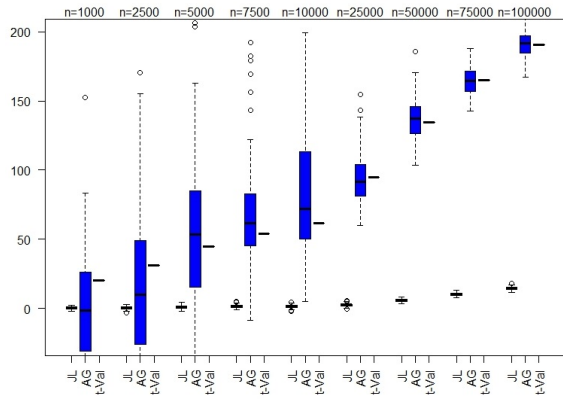
$$\begin{aligned}
 \sigma^2 &\leq \left(\frac{1}{\sqrt{n-p} - 2\sqrt{\ln(4/\alpha)}}\right)^2 \\
 \left(\|\widehat{\zeta}\|^2 + O\left(\left(1 + \frac{B^2 \sqrt{p+B}}{1-\eta}\right) \Delta \sqrt{\ln(1/\nu)} + \Delta^2 \|\widehat{X}^\top \widehat{X}^{-1}\|_F \cdot \ln(p/\nu)\right)\right)
 \end{aligned}$$

and indeed, the RHS is the definition of ρ^2 in the statement of Theorem D.1. \square

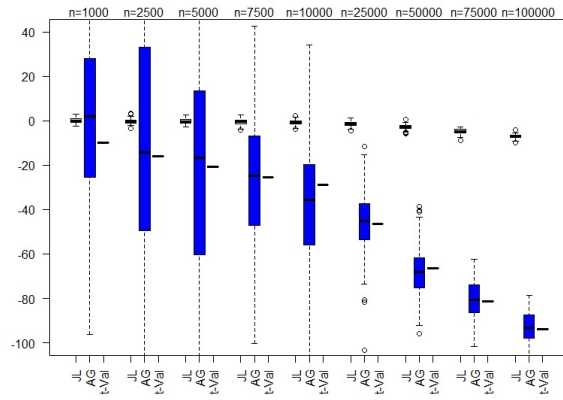
E. Experiment: Additional Figures

To complete our discussion about the experiments we have conducted, we attach here additional figures, plotting both the t -value approximations we get from both algorithms, and the ‘‘high-level decision’’ of whether correctly reject or not-reject the null hypothesis (and with what sign). First, we show the distribution of the t -value approximation for coordinates that should be rejected, in Figure 2, and then the decision of whether to reject or not based on this t -value — and whether it was right, conservative (we didn't reject while we needed to) or wrong (we rejected with the wrong sign, or rejected when we shouldn't have rejected) in Figure 3. As one can see, Algorithm 1 has far lower t -values (as expected) and therefore is much more conservative. In fact, it tends to not-reject coordinate 1 of the real-data even on the largest value of n (Figure 3c).

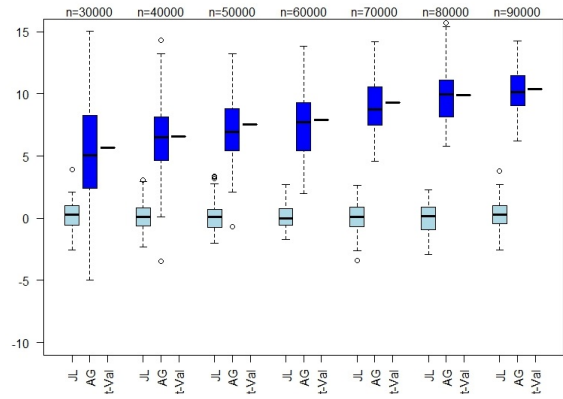
However, because Algorithm 1 also has much smaller variance, it also does not reject when it ought to not-reject, whereas Algorithm 2 erroneously rejects the null-hypotheses. This can be seen in Figures 4 and 5.



(a) Synthetic data, coordinate $\beta_1 = 0.5$

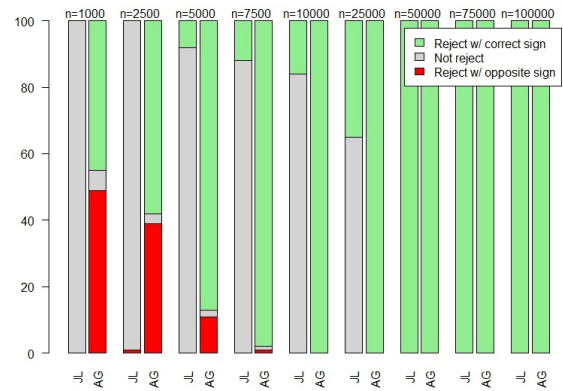


(b) Synthetic data, coordinate $\beta_2 = -0.25$

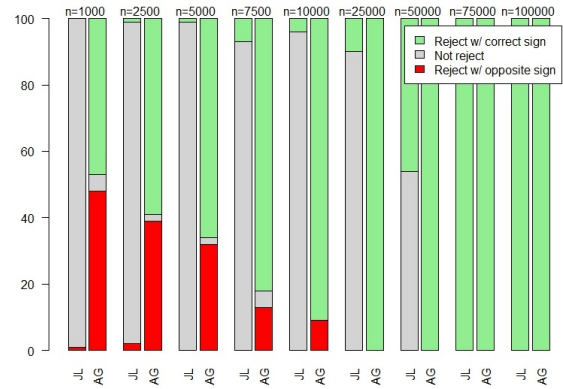


(c) real-life data, coordinate $\beta_1 = 14.07$

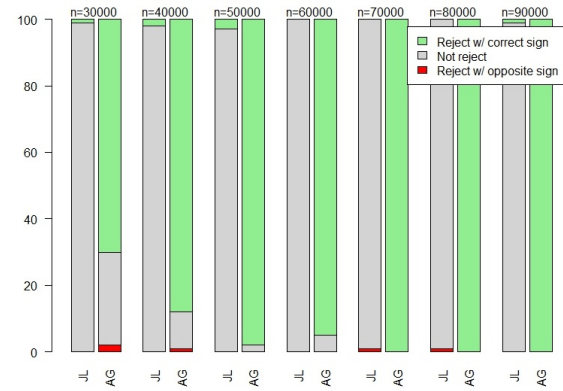
Figure 2. The distribution of the t -value approximations from selected experiments on synthetic and real-life data where the null hypothesis should be rejected



(a) Synthetic data, coordinate $\beta_1 = 0.5$

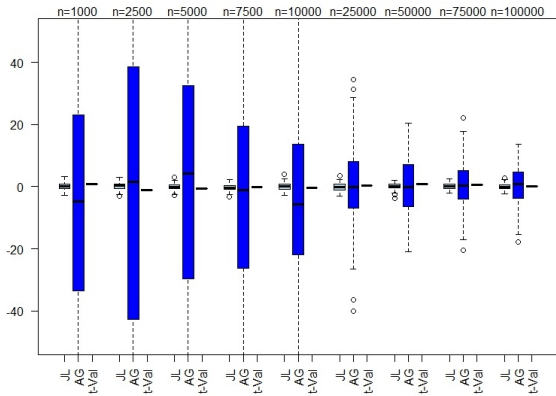


(b) Synthetic data, coordinate $\beta_2 = -0.25$

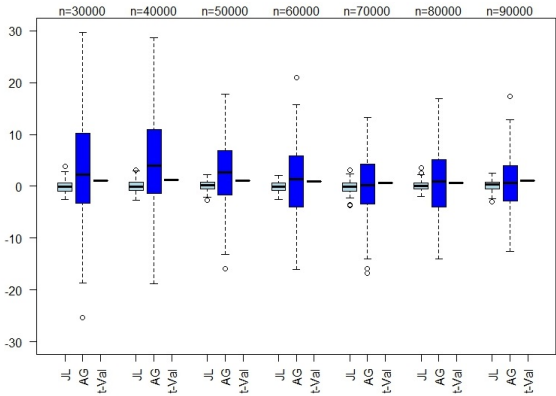


(c) real-life data, coordinate $\beta_1 = 14.07$

Figure 3. The correctness of our decision to reject the null-hypothesis based on the approximated t -value where the null hypothesis should be rejected

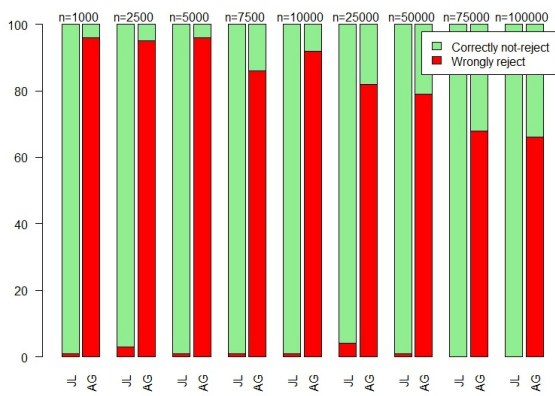


(a) Synthetic data, coordinate $\beta_3 = 0$

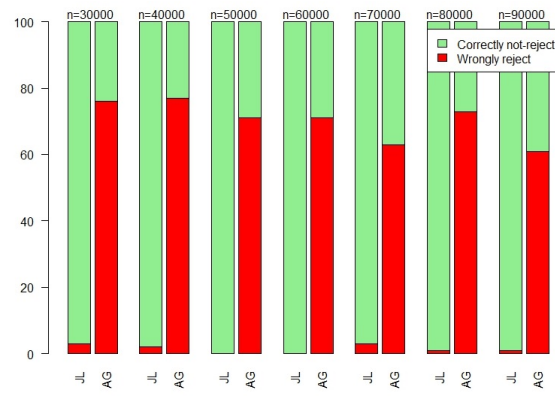


(b) real-life data, coordinate $\beta_2 = 0.57$

Figure 4. The distribution of the t -value approximations from selected experiments on synthetic and real-life data when the null hypothesis is (essentially) true



(a) Synthetic data, coordinate $\beta_3 = 0$



(b) Synthetic data, coordinate $\beta_2 = 0.57$

Figure 5. The correctness of our decision to reject the null hypothesis based on the approximated t -value when the null hypothesis is (essentially) true