

A PROOFS OF MAIN THEOREMS

A.1 Proof of Theorem 1

We fix $t = t(n, \boldsymbol{\lambda}) > 0$ to a positive value depending on $\boldsymbol{\lambda}$ and n that will be determined later. We define the following events:

1. For $i \in [n]$, \mathcal{B}_i is the event that $\mathbf{K}_{\setminus i}$ is non-singular and

$$\mathbf{y}_i \mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \boldsymbol{\Lambda} \mathbf{z}_i \geq 1.$$

2. For $i \in [n]$, \mathcal{S}_i is the event that $\mathbf{K}_{\setminus i}$ is singular.

3. \mathcal{S} is the event that \mathbf{K} is singular.

4. $\mathcal{B} := \mathcal{S} \cup \bigcup_{i=1}^n (\mathcal{B}_i \cup \mathcal{S}_i)$.

Additionally, we define the event $\mathcal{E}_i(t)$, for every $i \in [n]$ and a given $t > 0$, that $\mathbf{K}_{\setminus i}$ is non-singular and

$$\left\| \boldsymbol{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i} \right\|_2^2 \geq \frac{1}{t}.$$

Note that if the event \mathcal{B} does not occur, then $\mathbf{Z} \boldsymbol{\Lambda} \mathbf{Z}^\top$ is non-singular, each $\mathbf{K}_{\setminus i}$ is non-singular, and

$$\mathbf{y}_i \mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \boldsymbol{\Lambda} \mathbf{z}_i < 1, \quad \text{for all } i = 1, \dots, n.$$

Hence, by Lemma 1, if \mathcal{B} does not occur, then every training example is a support vector.

So, it suffices to upper-bound the probability of the event \mathcal{B} . We bound $\Pr(\mathcal{B})$ as follows:

$$\begin{aligned} \Pr(\mathcal{B}) &\leq \Pr(\mathcal{S}) + \sum_{i=1}^n \Pr(\mathcal{B}_i \cup \mathcal{S}_i) \\ &= \Pr(\mathcal{S}) + \sum_{i=1}^n \left(\Pr((\mathcal{B}_i \cap \mathcal{S}_i^c \cap \mathcal{E}_i(t)^c) \cup (\mathcal{S}_i \cap \mathcal{E}_i(t)^c)) + \Pr((\mathcal{B}_i \cup \mathcal{S}_i) \cap \mathcal{E}_i(t)) \right) \\ &\leq \Pr(\mathcal{S}) + \sum_{i=1}^n \left(\Pr(\mathcal{B}_i \mid \mathcal{S}_i^c \cap \mathcal{E}_i(t)^c) \Pr(\mathcal{S}_i^c \cap \mathcal{E}_i(t)^c) + \Pr(\mathcal{S}_i \cap \mathcal{E}_i(t)^c) + \Pr((\mathcal{B}_i \cup \mathcal{S}_i) \cap \mathcal{E}_i(t)) \right) \\ &\leq \Pr(\mathcal{S}) + \sum_{i=1}^n \left(\Pr(\mathcal{B}_i \mid \mathcal{S}_i^c \cap \mathcal{E}_i(t)^c) + \Pr(\mathcal{S}_i) + \Pr(\mathcal{E}_i(t)) \right). \end{aligned} \tag{7}$$

Above, the first two inequalities follow from the union bound, and the rest uses the law of total probability.

We first upper bound the probability of the singularity events in the following lemma.

Lemma 2. *We have*

$$\max\{\Pr(\mathcal{S}), \Pr(\mathcal{S}_1), \dots, \Pr(\mathcal{S}_n)\} \leq 2 \cdot 9^n \cdot \exp\left(-c \cdot \min\left\{\frac{d_2}{v^2}, \frac{d_\infty}{v}\right\}\right)$$

where $c > 0$ is the universal constant in the statement of Lemma 8.

Proof. It suffices to bound $\Pr(\mathcal{S})$, since each $\mathbf{K}_{\setminus i}$ is a principal submatrix of \mathbf{K} , and hence $\lambda_{\min}(\mathbf{K}_{\setminus i}) \geq \lambda_{\min}(\mathbf{K})$ for all $i \in [n]$. Observe that

$$\mathbf{Z} \boldsymbol{\Lambda} \mathbf{Z}^\top = \sum_{j=1}^d \lambda_j \mathbf{v}_j \mathbf{v}_j^\top$$

where \mathbf{v}_j is the j^{th} column of \mathbf{Z} . Recall that the columns of \mathbf{Z} are independent, and so these vectors satisfy the conditions of Lemma 8. Moreover, since $\mathbf{Z}\mathbf{\Lambda}\mathbf{Z}^\top$ is positive semi-definite, its singularity would require

$$\|\mathbf{Z}\mathbf{\Lambda}\mathbf{Z}^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}\|_2 \geq \|\boldsymbol{\lambda}\|_1.$$

The probability of this latter event can be bounded by Lemma 8 with $\tau = \|\boldsymbol{\lambda}\|_1$, thereby giving the claimed bound on $\Pr(\mathcal{S})$. This completes the proof of the lemma. \square

The next lemma upper bounds the probability of the event \mathcal{B}_i conditioned on the non-singularity event \mathcal{S}_i and the complement of the event $\mathcal{E}_i(t)$.

Lemma 3. *For any $t > 0$,*

$$\Pr(\mathcal{B}_i \mid \mathcal{S}_i^c \cap \mathcal{E}_i(t)^c) \leq 2 \exp\left(-\frac{t}{2v}\right).$$

Proof. Let \mathcal{B}'_i be the event that $\mathbf{K}_{\setminus i}$ is non-singular and

$$|\mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \mathbf{\Lambda} \mathbf{z}_i| = \max\{-\mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \mathbf{\Lambda} \mathbf{z}_i, \mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \mathbf{\Lambda} \mathbf{z}_i\} \geq 1.$$

Since $|y_i| = 1$, it follows that $\mathcal{B}_i \subseteq \mathcal{B}'_i$, so

$$\Pr(\mathcal{B}_i \mid \mathcal{S}_i^c \cap \mathcal{E}_i(t)^c) \leq \Pr(\mathcal{B}'_i \mid \mathcal{S}_i^c \cap \mathcal{E}_i(t)^c).$$

Conditional on the event $\mathcal{S}_i^c \cap \mathcal{E}_i(t)^c$, we have that $\mathbf{K}_{\setminus i}$ is non-singular and $\|\mathbf{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 \leq 1/t$. Since \mathbf{z}_i is independent of $\{(\mathbf{z}_j, y_j) : j \neq i\}$, it follows that

$$\mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \mathbf{\Lambda} \mathbf{z}_i = (\mathbf{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i})^\top \mathbf{z}_i$$

is (conditionally) sub-Gaussian with parameter at most $v \cdot \|\mathbf{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 \leq v/t$. Then, the standard sub-Gaussian tail bound gives us

$$\Pr(\mathcal{B}_i \mid \mathcal{S}_i^c \cap \mathcal{E}_i(t)^c) \leq \Pr(\mathcal{B}'_i \mid \mathcal{S}_i^c \cap \mathcal{E}_i(t)^c) \leq 2 \exp\left(-\frac{t}{2v}\right).$$

This completes the proof of the lemma. \square

Finally, the following lemma upper bounds the probability of the event $\mathcal{E}_i(t)$ for $t := d_\infty/2n$.

Lemma 4.

$$\Pr(\mathcal{E}_i(d_\infty/(2n))) \leq 2 \cdot 9^{n-1} \cdot \exp\left(-c \cdot \min\left\{\frac{d_2}{4v^2}, \frac{d_\infty}{v}\right\}\right)$$

where $c > 0$ is the universal constant from Lemma 8.

Proof. Let $\mathcal{E}'_i(t)$ be the event that

$$\lambda_{\min}(\mathbf{K}_{\setminus i}) \leq n \|\boldsymbol{\lambda}\|_\infty t.$$

Under \mathcal{S}_i^c , the matrix $\mathbf{K}_{\setminus i}$ is non-singular. We get

$$\begin{aligned} \|\mathbf{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 &\leq \|\mathbf{\Lambda}^{1/2}\|_{\text{op}}^2 \|\mathbf{\Lambda}^{1/2} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 \\ &= \|\boldsymbol{\lambda}\|_\infty \mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \mathbf{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i} \\ &\leq n \|\boldsymbol{\lambda}\|_\infty \sup_{\mathbf{u} \in \mathbb{R}^{n-1}: \|\mathbf{u}\|_2=1} \mathbf{u}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{u} \\ &= \frac{n \|\boldsymbol{\lambda}\|_\infty}{\lambda_{\min}(\mathbf{K}_{\setminus i})}. \end{aligned}$$

It follows that $\mathcal{E}_i(t) \subseteq \mathcal{E}'_i(t)$. Observe that for $t := d_\infty/(2n)$, the event $\mathcal{E}'_i(t)$ is that where

$$\lambda_{\min}(\mathbf{K}_{\setminus i}) \leq \frac{1}{2} \|\boldsymbol{\lambda}\|_1.$$

Therefore (as in the proof of Lemma 2), Lemma 8 with $\tau = \|\boldsymbol{\lambda}\|_1/2$ implies that

$$\begin{aligned} \Pr(\mathcal{E}'_i(d_\infty/(2n))) &= \Pr\left(\lambda_{\min}(\mathbf{K}_{\setminus i}) \leq \frac{1}{2} \|\boldsymbol{\lambda}\|_1\right) \\ &\leq 2 \cdot 9^{n-1} \cdot \exp\left(-c \cdot \min\left\{\frac{d_2}{4v^2}, \frac{d_\infty}{v}\right\}\right). \end{aligned}$$

This completes the proof of the lemma. \square

Plugging the probability bounds from Lemma 2, Lemma 3 and Lemma 4 (with $t = d_\infty/(2n)$) into Eq. (7) completes the proof of Theorem 1. \square

A.2 Proof of Theorem 2

The proof follows a similar sequence of steps to that of Theorem 1 with slight differences in the events that we condition on. We first observe that $\frac{1}{\sqrt{d}}\mathbf{z}_i \mid (\mathbf{Z}_{\setminus i}, \mathbf{y}_{\setminus i})$ is a uniformly random unit vector in S^{d-1} restricted to the subspace orthogonal to the row space of $\mathbf{Z}_{\setminus i}$. That is, it has the same (conditional) distribution as $\mathbf{B}_i \mathbf{u}_i$, where:

1. \mathbf{B}_i is a $d \times (n - d + 1)$ matrix whose columns form an orthonormal basis for the orthogonal complement of $\mathbf{Z}_{\setminus i}$'s row space;
2. \mathbf{u}_i is a uniformly random unit vector in S^{d-n} .

As before, for every $i \in [n]$, we define the event \mathcal{B}_i that $\mathbf{K}_{\setminus i}$ is non-singular and

$$\mathbf{y}_i \mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \boldsymbol{\Lambda} \mathbf{z}_i \geq 1.$$

The Haar measure ensures that the matrices \mathbf{Z} and $\mathbf{Z}_{\setminus i}$ always have full row rank. Therefore, because $\boldsymbol{\Lambda} \succ \mathbf{0}$, the matrices \mathbf{K} and $\mathbf{K}_{\setminus i}$ are always non-singular. So we do not need to worry about singularity (c.f. the events \mathcal{S} and \mathcal{S}_i). We accordingly consider the event $\mathcal{B} := \bigcup_{i=1}^n \mathcal{B}_i$. As before, we also define the event $\mathcal{E}_i(t)$ for every $i \in [n]$ and a given $t > 0$, that

$$\|\mathbf{B}_i^\top \boldsymbol{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 \geq \frac{d - n + 1}{d} \cdot \frac{1}{t}.$$

By the union bound, we get

$$\begin{aligned} \Pr(\mathcal{B}) &\leq \sum_{i=1}^n \Pr(\mathcal{B}_i) \\ &\leq \sum_{i=1}^n \Pr(\mathcal{B}_i \mid \mathcal{E}_i(t)^c) + \Pr(\mathcal{E}_i(t)), \end{aligned}$$

and so we need to upper bound the probabilities $\Pr(\mathcal{B}_i \mid \mathcal{E}_i(t)^c)$ and $\Pr(\mathcal{E}_i(t))$ for every $i \in [n]$.

The following lemma upper bounds $\Pr(\mathcal{B}_i \mid \mathcal{E}_i(t)^c)$, and is analogous to Lemma 3 in the proof of Theorem 1.

Lemma 5. *For any $t > 0$, we have*

$$\Pr(\mathcal{B}_i \mid \mathcal{E}_i(t)^c) \leq 2 \exp(-t).$$

Proof. First, as discussed above, we have

$$\begin{aligned} \Pr\left(y_i \mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \boldsymbol{\Lambda} \mathbf{z}_i \geq 1\right) &= \Pr\left(\sqrt{d} \cdot y_i \mathbf{y}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \boldsymbol{\Lambda} \mathbf{B}_i \mathbf{u}_i \geq 1\right) \\ &\leq \Pr\left(\sqrt{d} |(\mathbf{B}_i^\top \boldsymbol{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i})^\top \mathbf{u}_i| \geq 1\right). \end{aligned}$$

Moreover, \mathbf{u}_i is independent of $\mathbf{Z}_{\setminus i}$, and as established in Lemma 9, the random vector \mathbf{u}_i is sub-Gaussian with parameter at most $\mathcal{O}(1/(d-n+1))$. Therefore, $\sqrt{d} \cdot (\mathbf{B}_i^\top \boldsymbol{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i})^\top \mathbf{u}_i$ is conditionally sub-Gaussian with parameter at most $\frac{d}{d-n+1} \cdot \|\mathbf{B}_i^\top \boldsymbol{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 \leq \frac{1}{t}$. Here, the last inequality follows because we have conditioned on $\mathcal{E}_i(t)^c$. Therefore, the standard sub-Gaussian tail bound gives us

$$\Pr(\mathcal{B}_i \mid \mathcal{E}_i(t)^c) \leq 2 \exp(-t). \quad \square$$

The next lemma upper bounds $\Pr(\mathcal{E}_i(t))$ for $t := \frac{d-n+1}{d} \cdot \frac{d_\infty}{2n}$, and is analogous to Lemma 4 in the proof of Theorem 1.

Lemma 6. *We have*

$$\Pr\left(\mathcal{E}_i\left(\frac{d-n+1}{d} \cdot \frac{d_\infty}{2n}\right)\right) \leq \exp(-c_1 \cdot d) + 2 \cdot 9^n \cdot \exp(-c_2 \cdot \min\{d_2, d_\infty\})$$

where $c_1 > 0$ and $c_2 > 0$ are universal constants.

Proof. We get

$$\begin{aligned} \|\mathbf{B}_i^\top \boldsymbol{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 &\leq \|\mathbf{B}_i^\top\|_2^2 \cdot \|\boldsymbol{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 \\ &= \|\boldsymbol{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 \\ &\leq \frac{n \|\boldsymbol{\lambda}\|_\infty}{\lambda_{\min}(\mathbf{K}_{\setminus i})}, \end{aligned}$$

where we used the fact that \mathbf{B}_i has orthonormal columns, and the last inequality follows by an identical argument to the proof of Lemma 4. We will show in particular that

$$\Pr\left(\lambda_{\min}(\mathbf{K}_{\setminus i}) \geq \frac{1}{2} \|\boldsymbol{\lambda}\|_1\right) \geq 1 - \exp(-c_1 \cdot d) - 2 \cdot 9^{n-1} \cdot \exp(-c_2 \cdot \min\{d_2, d_\infty\}). \quad (8)$$

Given Eq. (8), we can complete the proof of Lemma 6. This is because we get

$$\|\mathbf{B}_i^\top \boldsymbol{\Lambda} \mathbf{Z}_{\setminus i}^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{y}_{\setminus i}\|_2^2 \leq \frac{2n \|\boldsymbol{\lambda}\|_\infty}{\|\boldsymbol{\lambda}\|_1} = \frac{2n}{d_\infty} = \frac{d-n+1}{d} \cdot \frac{1}{t}$$

for

$$t := \frac{d-n+1}{d} \cdot \frac{d_\infty}{2n}.$$

We complete the proof by proving Eq. (8). Let $\mathbf{S} \in \mathbb{R}^{m \times d}$ be a random matrix with iid standard Gaussian entries with $m := n-1$, and let the singular value decomposition of \mathbf{S} be $\mathbf{S} = \mathbf{V} \boldsymbol{\Lambda}_S \mathbf{U}^\top$ where $\mathbf{V} \in \mathbb{R}^{m \times m}$ and $\mathbf{U} \in \mathbb{R}^{d \times m}$ are orthonormal matrices. Then, it is well-known that $\sqrt{d} \cdot \mathbf{U}^\top$ follows the same distribution as $\mathbf{Z}_{\setminus i}$, and hence $\lambda_{\min}(\mathbf{K}_{\setminus i})$ has the same distribution as $d \cdot \lambda_{\min}(\mathbf{U}^\top \boldsymbol{\Lambda} \mathbf{U})$. Moreover,

$$\begin{aligned} d \cdot \lambda_{\min}(\mathbf{U}^\top \boldsymbol{\Lambda} \mathbf{U}) &= \min_{\mathbf{v} \in \mathbb{R}^n, \|\mathbf{v}\|_2=1} \mathbf{v}^\top \boldsymbol{\Lambda}_S^{-1} \mathbf{V}^\top \mathbf{V} \boldsymbol{\Lambda}_S \mathbf{U}^\top \boldsymbol{\Lambda} \mathbf{U} \boldsymbol{\Lambda}_S \mathbf{V}^\top \mathbf{V} \boldsymbol{\Lambda}_S^{-1} \mathbf{v} \\ &\geq \frac{d}{\|\boldsymbol{\Lambda}_S\|_{\text{op}}^2} \cdot \min_{\mathbf{v} \in \mathbb{R}^n, \|\mathbf{v}\|_2=1} \mathbf{v}^\top \mathbf{S} \boldsymbol{\Lambda} \mathbf{S}^\top \mathbf{v} \\ &= \frac{d}{\|\boldsymbol{\Lambda}_S\|_{\text{op}}^2} \cdot \lambda_{\min}(\mathbf{S} \boldsymbol{\Lambda} \mathbf{S}^\top). \end{aligned}$$

By classical operator norm tail bounds on Gaussian random matrices (e.g., Vershynin, 2010, Corollary 5.35), we note that $\|\mathbf{A}_S\|_2^2 \leq \frac{3}{2}d$ with probability at least $1 - \exp(-c_1 \cdot d)$. Now, we note that the matrix $\mathbf{S}\mathbf{A}\mathbf{S}^\top := \sum_{j=1}^d \lambda_j \mathbf{s}_j \mathbf{s}_j^\top$ where the \mathbf{s}_j 's are iid standard Gaussian random vectors in \mathbb{R}^n . So, we directly substitute Lemma 8 with $\tau := \frac{1}{4}\|\boldsymbol{\lambda}\|_1$, and get $\lambda_{\min}(\mathbf{S}\mathbf{A}\mathbf{S}^\top) \geq \frac{3}{4}\|\boldsymbol{\lambda}\|_1$ with probability at least $1 - 2 \cdot 9^m \cdot \exp(-c_2 \cdot \min\{d_2, d_\infty\})$. Putting both of these inequalities together directly gives us Eq. (8) with the desired probability bound, and completes the proof. \square

Finally, putting the high probability statements of Lemma 5 and Lemma 6 together completes the proof of Theorem 2.

A.3 Proof of Theorem 3

By Lemma 1, our task is equivalent to lower-bounding the probability that there exists $i \in [n]$ such that $y_i \mathbf{y}_i^\top (\mathbf{Z}_{\setminus i} \mathbf{Z}_{\setminus i}^\top)^{-1} \mathbf{Z}_{\setminus i} \mathbf{z}_i \geq 1$. This event is the union of n (possibly overlapping) events, and hence its probability is at least the probability of one of the events, say, the first one:

$$\Pr\left(\exists i \in [n] \text{ s.t. } y_i \mathbf{y}_i^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \mathbf{z}_i \geq 1\right) \geq \Pr\left(y_1 \mathbf{y}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \mathbf{z}_1 \geq 1\right).$$

Because \mathbf{z}_1 is a standard Gaussian random vector independent of $\mathbf{Z}_{\setminus 1}$, the conditional distribution of $y_1 \mathbf{y}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \mathbf{z}_1 \mid \mathbf{Z}_{\setminus 1}$ is Gaussian with mean zero and variance $\sigma^2 := \|\mathbf{Z}_{\setminus 1}^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{y}_1\|_2^2$. Therefore, for any $t > 0$, we have

$$\begin{aligned} \Pr\left(y_1 \mathbf{y}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \mathbf{z}_1 \geq 1\right) &= \mathbb{E}\left[\Pr(\sigma g \geq 1 \mid \sigma)\right] \quad (\text{where } g \sim \mathcal{N}(0, 1), g \perp \sigma) \\ &= \mathbb{E}\left[\Phi(-1/\sigma)\right] \\ &\geq \mathbb{E}\left[\Phi(-1/\sigma) \mid \sigma^2 \geq 1/t\right] \Pr\left(\sigma^2 \geq 1/t\right) \\ &\geq \Phi(-\sqrt{t}) \cdot \Pr(\mathcal{E}_1(t)), \end{aligned}$$

where Φ is the standard Gaussian cumulative distribution function, and $\mathcal{E}_1(t)$ is the event that

$$\sigma^2 = \mathbf{y}_{\setminus 1} \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \mathbf{Z}_{\setminus 1}^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{y}_{\setminus 1} = \mathbf{y}_{\setminus 1} \mathbf{K}_{\setminus 1}^{-1} \mathbf{y}_{\setminus 1} \geq \frac{1}{t}$$

(as in the proofs of Theorem 1 and Theorem 2). We now lower-bound the probability of $\mathcal{E}_1(t)$. Observe that the $(n-1) \times (n-1)$ random matrix $\mathbf{K}_{\setminus 1} = \mathbf{Z}_{\setminus 1} \mathbf{Z}_{\setminus 1}^\top$ follows a Wishart distribution with identity scale matrix and d degrees-of-freedom. Moreover, by the rotational symmetry of the standard Gaussian distribution, the random variable $\mathbf{y}_{\setminus 1}^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{y}_{\setminus 1}$ has the same distribution as that of $(\sqrt{n-1} \mathbf{e}_1)^\top \mathbf{K}_{\setminus 1}^{-1} (\sqrt{n-1} \mathbf{e}_1) = (n-1) \mathbf{e}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{e}_1$. It is known that $1/\mathbf{e}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{e}_1$ follows a χ^2 distribution with $d - (n-2)$ degrees-of-freedom; we denote its cumulative distribution function by F_{d-n+2} . Therefore,

$$\Pr(\mathcal{E}_1(t)) = F_{d-n+2}(t(n-1)).$$

So, we have shown that

$$\Pr\left(y_1 \mathbf{y}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \mathbf{z}_1 \geq 1\right) \geq \sup_{t \geq 0} \Phi(-\sqrt{t}) \cdot F_{d-n+2}(t(n-1)).$$

For $t := \frac{d-n+4+2\sqrt{d-n+2}}{n-1}$, we obtain $F_{d-n+2}(t) \geq 1 - 1/e$ by a standard χ^2 tail bound (Laurent and Massart, 2000, Lemma 1). In this case, we obtain

$$\Pr\left(y_1 \mathbf{y}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \mathbf{z}_1 \geq 1\right) \geq \Phi\left(-\sqrt{\frac{d-n+4+2\sqrt{d-n+2}}{n-1}}\right) \cdot \left(1 - \frac{1}{e}\right) \quad (9)$$

as claimed. \square

B ANISOTROPIC VERSION OF THEOREM 3

Below, we give a version of Theorem 3 that applies to certain anisotropic settings, depending on some conditions on λ .

Theorem 4. *There are absolute constants $c > 0$ and $c' > 0$ such that the following hold. Let the training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ follow the model from Section 2.2, with $\mathbf{z}_1, \dots, \mathbf{z}_n$ being iid standard Gaussian random vectors in \mathbb{R}^d , and $y_1, \dots, y_n \in \{\pm 1\}$ being arbitrary but fixed (i.e., non-random) values. Assume $d > n$ and that there exists $k \in \mathbb{N}$ and $b > 1$ such that $k < (n-1)/c$ and*

$$\frac{\sum_{j=k+1}^d \lambda_j}{\lambda_{k+1}} \leq b(n-1)$$

where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$. Then the probability that at least one training example is not a support vector is at least

$$c' \cdot \Phi \left(-\sqrt{\frac{2cb^2(n-1)}{k+1}} \right) \cdot \left(1 - 10e^{-(n-1)/c} \right),$$

where Φ is the standard Gaussian cumulative distribution function.

Note that the probability bound in Theorem 4 is at least a positive constant for sufficiently large n provided that the (k, b) obtained as a function of λ satisfy $k+1 \geq c''b^2(n-1)$ for some absolute constant $c'' > 0$.

Proof. The proof begins in the same way as in that of Theorem 3. Using the same arguments, we obtain the following lower bound:

$$\begin{aligned} \Pr \left(\exists i \in [n] \text{ s.t. } y_i \mathbf{y}_i^\top \mathbf{K}_{\setminus i}^{-1} \mathbf{Z}_{\setminus i} \boldsymbol{\Lambda} \mathbf{z}_i \geq 1 \right) &\geq \Pr \left(y_1 \mathbf{y}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \boldsymbol{\Lambda} \mathbf{z}_1 \geq 1 \right) \\ &\geq \Phi(-\sqrt{t}) \cdot \Pr(\mathcal{E}_1(t)) \end{aligned} \quad (10)$$

where $\mathcal{E}_1(t)$ is the event that

$$\left\| \boldsymbol{\Lambda} \mathbf{Z}_{\setminus 1}^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{y}_{\setminus 1} \right\|_2^2 \geq \frac{1}{t}.$$

We next focus on lower-bounding the probability of $\mathcal{E}_1(t)$. (This part is more involved than in the proof of Theorem 3.) Observe that the (rotationally invariant) distribution of $\mathbf{Z}_{\setminus 1}$ is the same as that of $\mathbf{Q} \mathbf{Z}_{\setminus 1}$, where \mathbf{Q} is a uniformly random $(n-1) \times (n-1)$ orthogonal matrix independent of $\mathbf{Z}_{\setminus 1}$. Therefore, $\boldsymbol{\Lambda} \mathbf{Z}_{\setminus 1}^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{y}_{\setminus 1}$ has the same distribution as

$$\begin{aligned} \boldsymbol{\Lambda} (\mathbf{Q} \mathbf{Z}_{\setminus 1})^\top (\mathbf{Q} \mathbf{Z}_{\setminus 1} \boldsymbol{\Lambda} \mathbf{Z}_{\setminus 1}^\top \mathbf{Q}^\top)^{-1} \mathbf{y}_{\setminus 1} &= \boldsymbol{\Lambda} \mathbf{Z}_{\setminus 1}^\top \mathbf{Q}^\top \mathbf{Q} (\mathbf{Z}_{\setminus 1} \boldsymbol{\Lambda} \mathbf{Z}_{\setminus 1}^\top)^{-1} \mathbf{Q}^\top \mathbf{y}_{\setminus 1} \\ &= \sqrt{n-1} \boldsymbol{\Lambda} \mathbf{Z}_{\setminus 1}^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{u} \end{aligned}$$

where $\mathbf{u} := \mathbf{Q}^\top \mathbf{y}_{\setminus 1} / \sqrt{n-1}$ is a uniformly random unit vector, independent of $\mathbf{Z}_{\setminus 1}$. Letting $\mathbf{M} := \boldsymbol{\Lambda} \mathbf{Z}_{\setminus 1}^\top \mathbf{K}_{\setminus 1}^{-1}$, we can thus lower-bound the probability of $\mathcal{E}_1(t)$ using

$$\begin{aligned} \Pr(\mathcal{E}_1(t)) &= \Pr \left(\left\| \sqrt{n-1} \mathbf{M} \mathbf{u} \right\|_2^2 > 1/t \right) \\ &\geq \Pr \left(\left\| \sqrt{n-1} \mathbf{M} \mathbf{u} \right\|_2^2 > 1/t \mid \text{tr}(\mathbf{M}^\top \mathbf{M}) \geq 2/t \right) \cdot \Pr \left(\text{tr}(\mathbf{M}^\top \mathbf{M}) \geq 2/t \right). \end{aligned} \quad (11)$$

We lower-bound each of the probabilities on the right-hand side of Eq. (11).

We begin with the first probability in Eq. (11), which we handle for arbitrary $t > 0$. By the Paley-Zygmund inequality, we have

$$\Pr \left(\left\| \sqrt{n-1} \mathbf{M} \mathbf{u} \right\|_2^2 > \frac{1}{2} \mathbb{E} \left[\left\| \sqrt{n-1} \mathbf{M} \mathbf{u} \right\|_2^2 \right] \mid \mathbf{Z}_{\setminus 1} \right) \geq \frac{1}{4} \cdot \frac{\mathbb{E} \left[\left\| \sqrt{n-1} \mathbf{M} \mathbf{u} \right\|_2^2 \right]^2}{\mathbb{E} \left[\left\| \sqrt{n-1} \mathbf{M} \mathbf{u} \right\|_2^4 \right]}. \quad (12)$$

Since $\sqrt{n-1}\mathbf{u}$ is isotropic, we have

$$\mathbb{E} \left[\|\sqrt{n-1}\mathbf{M}\mathbf{u}\|_2^2 \mid \mathbf{Z}_{\setminus 1} \right] = (n-1) \operatorname{tr}(\mathbf{M}^\top \mathbf{M} \mathbb{E}[\mathbf{u}\mathbf{u}^\top]) = \operatorname{tr}(\mathbf{M}^\top \mathbf{M}).$$

Furthermore, by Lemma 9,

$$\mathbb{E} \left[\|\sqrt{n-1}\mathbf{M}\mathbf{u}\|_2^4 \mid \mathbf{Z}_{\setminus 1} \right] \leq C \operatorname{tr}(\mathbf{M}^\top \mathbf{M})^2$$

for some universal constant $C > 0$. Therefore, plugging back into Eq. (12), we obtain

$$\Pr \left(\|\sqrt{n-1}\mathbf{M}\mathbf{u}\|_2^2 > \frac{1}{2} \operatorname{tr}(\mathbf{M}^\top \mathbf{M}) \mid \mathbf{Z}_{\setminus 1} \right) \geq \frac{1}{4} \cdot \frac{\operatorname{tr}(\mathbf{M}^\top \mathbf{M})^2}{C \operatorname{tr}(\mathbf{M}^\top \mathbf{M})^2} = \frac{1}{4C}.$$

Thus we also have the following for arbitrary $t > 0$:

$$\Pr \left(\|\sqrt{n-1}\mathbf{M}\mathbf{u}\|_2^2 > 1/t \mid \operatorname{tr}(\mathbf{M}^\top \mathbf{M}) \geq 2/t \right) \geq \frac{1}{4C}. \quad (13)$$

We next consider the second probability in Eq. (11), namely $\Pr(\operatorname{tr}(\mathbf{M}^\top \mathbf{M}) \geq 2/t)$. Recall that we assume there exists $k < (n-1)/c$ and $b > 1$ such that

$$\frac{\sum_{j=k+1}^d \lambda_j}{\lambda_{k+1}} \leq b(n-1). \quad (14)$$

We claim that for $t := \frac{2cb^2(n-1)}{k+1}$,

$$\Pr \left(\operatorname{tr}(\mathbf{M}^\top \mathbf{M}) \geq \frac{2}{t} \right) \geq 1 - 10e^{-(n-1)/c}. \quad (15)$$

Indeed, this claim follows from Lemma 16 of (Bartlett et al., 2020), where their matrix \mathbf{C} is our matrix $\mathbf{M}^\top \mathbf{M}$, except our matrix is $(n-1) \times (n-1)$ instead of $n \times n$, and their matrix $\mathbf{\Sigma}$ is our matrix $\mathbf{\Lambda}$; see the definitions in their Lemma 8. The universal constant $c > 0$ in their lemma is the same as ours, and Eq. (14) is precisely their condition $r_k(\mathbf{\Sigma}) < b(n-1)$ (with the same k and b). Therefore, the conclusion of their lemma implies, in our notation, that with probability at least $1 - 10e^{-(n-1)/c}$,

$$\operatorname{tr}(\mathbf{M}^\top \mathbf{M}) \geq \frac{k+1}{cb^2(n-1)} = \frac{2}{t}.$$

This proves the claimed probability bound.

We conclude from Eq. (10), Eq. (11), Eq. (13), and Eq. (15), that the probability that at least one training example is not a support vector is bounded below by

$$\Phi \left(-\sqrt{\frac{2cb^2(n-1)}{k+1}} \right) \cdot \frac{1}{4C} \cdot \left(1 - 10e^{-(n-1)/c} \right)$$

as claimed. \square

C TIGHTNESS OF ARGUMENT IN THEOREM 3

We show below that our bound on $\Pr(y_1 \mathbf{y}_{\setminus 1}^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \mathbf{z}_1 \geq 1)$ from the proof of Theorem 3 is essentially tight. This means that in order to improve our converse result, we cannot only improve our bound on the aforementioned probability. It seems important to be able to handle simultaneously the conditions corresponding to multiple training examples, which our present arguments do not do. In particular, resolving this gap would require reasoning about whether the indicator random variables, that the conditions are violated, are highly correlated or not. If they are, we should expect the phase transition to happen at $d \sim n$ (as predicted by the converse); if they are not, we should expect the phase transition to happen at $d \sim n \log n$ (as predicted by the upper bound).

Carrying over the notation from the proof above, we have the following upper-bound:

$$\begin{aligned} \Pr\left(\mathbf{y}_1 \mathbf{y}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \mathbf{z}_1 \geq 1\right) &= \mathbb{E}\left(\Phi(-1/\sigma)\right) \\ &\leq \inf_{t>0} \left\{ \Phi(-\sqrt{t}) + \Pr(\mathcal{E}_1(t)) \right\}. \end{aligned}$$

The last step follows by the law of total probability, and noting that $\Phi(-x)$ is a decreasing function in x as well as being bounded above by one. We will bound the second term for a suitable choice of t . Recall that $\mathcal{E}_1(t)$ is the event that

$$\sigma^2 = \mathbf{y}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{y}_1 \geq \frac{1}{t}.$$

Observe that $\sigma^2 \leq \frac{n-1}{\lambda_{\min}(\mathbf{K}_{\setminus 1})}$, where the $(n-1) \times (n-1)$ random matrix $\mathbf{K}_{\setminus 1} = \mathbf{Z}_{\setminus 1} \mathbf{Z}_{\setminus 1}^\top$ follows a Wishart distribution with identity scale matrix and d degrees of freedom. Directly quoting (Vershynin, 2018, Theorem 5.32), we get

$$\Pr\left(\lambda_{\min}(\mathbf{K}_{\setminus 1}) \leq (\sqrt{d} - \sqrt{n} - \delta)^2\right) \leq e^{-\delta^2/2}.$$

for any value of δ such that $0 < \delta < \sqrt{d} - \sqrt{n}$. Therefore,

$$\Pr\left(\sigma^2 \geq \frac{n-1}{(\sqrt{d} - \sqrt{n} - \delta)^2}\right) \leq \Pr\left(\lambda_{\min}(\mathbf{K}_{\setminus 1}) \leq (\sqrt{d} - \sqrt{n} - \delta)^2\right) \leq e^{-\delta^2/2}.$$

Assuming $d > 4n$, we set $\delta := \sqrt{n}$ and $t := \frac{(\sqrt{d} - 2\sqrt{n})^2}{n-1}$, and obtain

$$\begin{aligned} \Pr\left(\mathbf{y}_1 \mathbf{y}_1^\top \mathbf{K}_{\setminus 1}^{-1} \mathbf{Z}_{\setminus 1} \mathbf{z}_1 \geq 1\right) &\leq \Phi(-\sqrt{t}) + \Pr(\mathcal{E}_1(t)) \\ &\leq \Phi\left(-\frac{\sqrt{d} - 2\sqrt{n}}{\sqrt{n-1}}\right) + e^{-n/2}, \end{aligned}$$

which can be directly compared to Eq. (9).

D PROBABILISTIC INEQUALITIES

Lemma 7. *Let $\mathbf{M} \in \mathbb{R}^{n \times n}$ be a symmetric matrix, and let \mathcal{N} be an ϵ -net of S^{n-1} with respect to the Euclidean metric for some $\epsilon < 1/2$. Then*

$$\|\mathbf{M}\|_2 \leq \frac{1}{1-2\epsilon} \max_{\mathbf{u} \in \mathcal{N}} |\mathbf{u}^\top \mathbf{M} \mathbf{u}|.$$

Proof. See (Vershynin, 2010, Lemma 5.4). □

Lemma 8. *There is a universal constant $c > 0$ such that the following holds. Let $\lambda_1, \dots, \lambda_d > 0$ be given. Let $\mathbf{v}_1, \dots, \mathbf{v}_d$ be independent random vectors taking values in \mathbb{R}^n such that, for some $v > 0$,*

$$\mathbb{E}(\mathbf{v}_j) = \mathbf{0}, \quad \mathbb{E}(\mathbf{v}_j \mathbf{v}_j^\top) = \mathbf{I}_n, \quad \mathbb{E}(\exp(\mathbf{u}^\top \mathbf{v}_j)) \leq \exp(v \|\mathbf{u}\|_2^2 / 2) \quad \text{for all } \mathbf{u} \in \mathbb{R}^n$$

for all $j = 1, \dots, d$. For any $\tau > 0$,

$$\Pr\left(\left\| \sum_{j=1}^d \lambda_j \mathbf{v}_j \mathbf{v}_j^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}_n \right\|_2 \geq \tau\right) \leq 2 \cdot 9^n \cdot \exp\left(-c \cdot \min\left\{\frac{\tau^2}{v^2 \|\boldsymbol{\lambda}\|_2^2}, \frac{\tau}{v \|\boldsymbol{\lambda}\|_\infty}\right\}\right).$$

where $\|\boldsymbol{\lambda}\|_1 := \sum_{j=1}^d \lambda_j$, $\|\boldsymbol{\lambda}\|_2^2 := \sum_{j=1}^d \lambda_j^2$, and $\|\boldsymbol{\lambda}\|_\infty := \max_{j \in [d]} \lambda_j$.

Proof. Let \mathcal{N} be an $(1/4)$ -net of S^{n-1} with respect to the Euclidean metric. A standard volume argument of Pisier (1999) allows a choice of \mathcal{N} with $|\mathcal{N}| \leq 9^n$. By Lemma 7, we have for any $t > 0$,

$$\Pr \left(\left\| \sum_{j=1}^d \lambda_j \mathbf{v}_j \mathbf{v}_j^\top - \|\boldsymbol{\lambda}\|_1 \mathbf{I}_n \right\|_2 \geq \tau \right) \leq \Pr \left(\max_{\mathbf{u} \in \mathcal{N}} \left| \sum_{j=1}^d \lambda_j (\mathbf{u}^\top \mathbf{v}_j)^2 - \|\boldsymbol{\lambda}\|_1 \right| \geq \tau/2 \right).$$

Next, observe that for any $\mathbf{u} \in S^{n-1}$, the random variables $\mathbf{u}^\top \mathbf{v}_1, \dots, \mathbf{u}^\top \mathbf{v}_d$ are independent random variables, each with mean-zero, unit variance, and sub-Gaussian with parameter v . By the Hanson-Wright inequality of Rudelson and Vershynin (2013) and a union bound, there exists a universal constant $c > 0$ such that, for any unit vector $\mathbf{u} \in S^{n-1}$ and any $\tau > 0$,

$$\Pr \left(\max_{\mathbf{u} \in \mathcal{N}} \left| \sum_{j=1}^d \lambda_j (\mathbf{u}^\top \mathbf{v}_j)^2 - \|\boldsymbol{\lambda}\|_1 \right| \geq \tau/2 \right) \leq 2 \cdot 9^n \cdot \exp \left(-c \cdot \min \left\{ \frac{\tau^2}{v^2 \|\boldsymbol{\lambda}\|_2^2}, \frac{\tau}{v \|\boldsymbol{\lambda}\|_\infty} \right\} \right).$$

The claim follows. \square

Lemma 9. *Let $\boldsymbol{\theta}$ be a uniformly random unit vector in S^{m-1} . For any unit vector $\mathbf{u} \in S^{m-1}$, the random variable $\mathbf{u}^\top \boldsymbol{\theta}$ is sub-Gaussian with parameter $v = O(1/m)$. Moreover, for any matrix $\mathbf{M} \in \mathbb{R}^{m \times m}$, we have*

$$\mathbb{E} \left[\|\mathbf{M}\boldsymbol{\theta}\|_2^4 \right] \leq \frac{C}{m^2} \text{tr}(\mathbf{M}^\top \mathbf{M})^2$$

where $C > 0$ is a universal constant.

Proof. Let L be a χ random variable with m degrees-of-freedom, independent of $\boldsymbol{\theta}$, so the distribution of $\mathbf{z} := L\boldsymbol{\theta}$ is the standard Gaussian in \mathbb{R}^m . Let $\mu := \mathbb{E}[L] = \mathbb{E}[L \mid \boldsymbol{\theta}] = \sqrt{2} \frac{\Gamma((m+1)/2)}{\Gamma(m/2)} = \Omega(\sqrt{m})$. By Jensen's inequality, for any $t \in \mathbb{R}$,

$$\begin{aligned} \mathbb{E} \left[\exp(t\mathbf{u}^\top \boldsymbol{\theta}) \right] &= \mathbb{E} \left[\exp \left(\left(\frac{t}{\mu} \mathbf{u} \right)^\top (\mathbb{E}[L \mid \boldsymbol{\theta} \boldsymbol{\theta}]) \right) \right] \\ &\leq \mathbb{E} \left[\exp \left(\left(\frac{t}{\mu} \mathbf{u} \right)^\top (L\boldsymbol{\theta}) \right) \right] \\ &= \mathbb{E} \left[\exp \left(\left(\frac{t}{\mu} \mathbf{u} \right)^\top \mathbf{z} \right) \right] \\ &= \exp \left(\frac{t^2}{2\mu^2} \right). \end{aligned}$$

It follows that $\mathbf{u}^\top \boldsymbol{\theta}$ is sub-Gaussian with parameter $v = 1/\mu^2 = O(1/m)$.

Similarly, again by Jensen's inequality,

$$\begin{aligned} \mu^4 \cdot \mathbb{E} \left[\|\mathbf{M}\boldsymbol{\theta}\|_2^4 \right] &= \mathbb{E} \left[\mathbb{E}[L \mid \boldsymbol{\theta}]^4 \|\mathbf{M}\boldsymbol{\theta}\|_2^4 \right] \\ &\leq \mathbb{E} \left[L^4 \|\mathbf{M}\boldsymbol{\theta}\|_2^4 \right] \\ &= \mathbb{E} \left[\|\mathbf{M}\mathbf{z}\|_2^4 \right]. \end{aligned}$$

Furthermore, a direct computation shows that

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{M}\mathbf{z}\|_2^4 \right] &= 2 \text{tr}((\mathbf{M}^\top \mathbf{M})^2) + \text{tr}(\mathbf{M}^\top \mathbf{M})^2 \\ &\leq 3 \text{tr}(\mathbf{M}^\top \mathbf{M})^2. \end{aligned}$$

The conclusion follows since $\mu^4 = \Omega(m^2)$. \square