

A Learning LSFs with Bounded Noise in Kendall's Tau distance

A.1 Improperly Learning LSFs with Bounded Noise

We provide an improper learner for LSFs in the presence of bounded noise. We first restate the main result of this section, whose proof relies on a connection between noisy linear label ranking distributions and the Massart noise model.

Theorem 3 (Non-Proper Learning Algorithm). *Fix $\eta \in [0, 1/2)$ and $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2. `ImproperLSF` (Algorithm 1) draws $N = \tilde{O}\left(\frac{d}{\epsilon(1-2\eta)^6} \log(k/\delta)\right)$ samples from \mathcal{D} , runs in $\text{poly}(d, k, 1/\epsilon, \log(1/\delta))$ time and, with probability at least $1 - \delta$, outputs a hypothesis $h : \mathbb{R}^d \rightarrow \mathbb{S}_k$ that is ϵ -close in KT distance to the target.*

Proof. Assume that the target function is $\sigma^*(\mathbf{x}) = \sigma_{\mathbf{W}^*}(\mathbf{x}) = \text{argsort}(\mathbf{W}^* \mathbf{x})$ for some unknown matrix $\mathbf{W}^* \in \mathbb{R}^{k \times d}$. Consider a collection of N i.i.d. samples from an η -noisy linear label ranking distribution \mathcal{D} (see Definition 2) and let T be the associated training set. For each example $(\mathbf{x}, \pi) \in T$, we create a list of $\binom{k}{2}$ binary examples (\mathbf{x}, y_{ij}) with $y_{ij} = \text{sgn}(\pi(i) - \pi(j))$ for any $1 \leq i < j \leq k$, where $\pi(i)$ denotes the position of the element i . Hence, we create the datasets T_{ij} consisting of the binary labeled examples (\mathbf{x}, y_{ij}) . We have that

$$\Pr_{(\mathbf{x}, \pi) \sim \mathcal{D}} [y_{ij} \cdot \text{sgn}((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x}) < 0 \mid \mathbf{x}] = \Pr_{\pi \sim \mathcal{M}(\sigma^*(\mathbf{x}))} [\pi(i) < \pi(j) \mid \mathbf{W}_i^* \cdot \mathbf{x} < \mathbf{W}_j^* \cdot \mathbf{x}].$$

Since $\mathcal{M}(\sigma^*(\mathbf{x}))$ is an η -bounded noise ranking distribution (see Definition 1), we get that

$$\Pr_{\pi \sim \mathcal{M}(\sigma^*(\mathbf{x}))} [\pi(i) < \pi(j) \mid \sigma^*(\mathbf{x})(i) > \sigma^*(\mathbf{x})(j)] \leq \eta < 1/2,$$

where $\sigma^*(\mathbf{x})(i)$ denotes the position of the element i in the ranking $\sigma^*(\mathbf{x})$. Focusing on the training set T_{ij} , we have that the sign y_{ij} is flipped with probability at most η . So, we have reduced the problem to $\binom{k}{2}$ sub-problems concerning the learnability of halfspaces in the presence of Massart noise. The Massart noise model is a special case of Definition 2 where $k = 2$. Note also that for each training set T_{ij} , the features \mathbf{x} have the same distribution. We can now apply the following result for LTFs with Massart noise for the standard Gaussian distribution. Recall that the concept class of homogeneous halfspaces (or linear threshold functions) is $\mathcal{C}_{\text{LTF}} = \{h_{\mathbf{w}}(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x}) : \mathbf{w} \in \mathbb{R}^d\}$.

Lemma 6 (Learning Halfspaces with Massart noise [ZSA20]). *Fix $\eta \in [0, 1/2)$ and let $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2 with $k = 2$ (where $\mathcal{C}_{\text{LSF}} = \mathcal{C}_{\text{LTF}}$). There is a computationally efficient algorithm `MassartLTF` that draws $m = O\left(\frac{d \cdot \text{polylog}(d)}{\epsilon(1-2\eta)^6} \cdot \log(1/\delta)\right)$ samples from \mathcal{D} , runs in $\text{poly}(m)$ time and outputs a linear threshold function h that is ϵ -close to the target linear threshold function h^* with probability at least $1 - \delta$, i.e., it holds $\Pr_{\mathbf{x} \sim \mathcal{N}_d} [h(\mathbf{x}) \neq h^*(\mathbf{x})] \leq \epsilon$.*

We can invoke the algorithm of Lemma 6 for any alternatives $1 \leq i < j \leq k$ with accuracy $\epsilon' = O(\epsilon)$, $\delta' = O(\delta/k^2)$ and error rate $\eta < 1/2^4$. We remark that Lemma 6 returns a halfspace. Each one of the $\binom{k}{2}$ calls will provide a vector $\mathbf{v}_{ij} \in \mathbb{R}^d$ such that, with probability at least $1 - \delta'$, it satisfies

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\text{sgn}(\mathbf{v}_{ij} \cdot \mathbf{x}) \neq \text{sgn}((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x})] \leq \epsilon',$$

where the true target halfspace has normal vector $\mathbf{W}_i^* - \mathbf{W}_j^*$. Moreover, for any $i < j$, the algorithm requires that the training set T_{ij} is of size

$$|T_{ij}| = \Omega\left(\frac{d}{\epsilon'} \cdot \frac{1}{(1-2\eta)^6} \cdot \log(1/\delta')\right),$$

and, so, a total number of

$$N = \Omega\left(\frac{d}{\epsilon} \cdot \frac{1}{(1-2\eta)^6} \cdot \log(k/\delta)\right),$$

⁴We can assume that η is known without loss of generality.

samples (\mathbf{x}, π) is required from the distribution \mathcal{D} . Given a collection of linear classifiers with normal vectors \mathbf{v}_{ij} for any $i < j$, it remains to aggregate them and compute a sorting function $h : \mathbb{R}^d \rightarrow \mathbb{S}_k$. To this end, the estimator h , given an example \mathbf{x} , creates the directed complete graph G with k nodes with directed edge $i \rightarrow j$ if $\mathbf{v}_{ij} \cdot \mathbf{x} > 0$. If all the linear classifiers are correct (which occurs with probability $1 - O(\epsilon k^2)$ over \mathcal{D}_x due to the union bound), the graph G is acyclic (since it will match the true directions induced by \mathbf{W}^*) and the estimator h outputs the induced permutation. Observe that the KT distance is

$$\frac{1}{\binom{k}{2}} \cdot \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} \left[\sum_{1 \leq i < j \leq k} \mathbf{1}\{\text{sgn}(\mathbf{v}_{ij} \cdot \mathbf{x}) \neq \text{sgn}((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x})\} \right] \leq \epsilon'.$$

Otherwise, the classifiers are inconsistent and G contains cycles. So, the expected number of mistakes in the graph G is ϵk^2 . The estimator in order to output a ranking uses a deterministic constant approximation algorithm for the minimum Feedback Arc Set [ACN08] in order to remove the cycles. For an overview of this fundamental line of research, we refer to [ACN08, VZW09, KMS06].

Lemma 7 (3-Approximation Algorithm for minimum FAS (see [VZW09, ACN08])). *There is a deterministic algorithm MFAS for the minimum Feedback Arc Set on unweighted tournaments with k vertices that outputs orderings with cost less than $3 \cdot \text{OPT}$. The running time is $\text{poly}(k)$.*

In the above, OPT is the minimum number of flips the algorithm should perform. With input the cyclic directed graph G induced by the estimated linear classifiers, the algorithm of Lemma 7 computes, in $\text{poly}(k)$ time, a 3-approximation of the optimal solution (i.e., instead of correcting ϵ_0 directed edges, the algorithm will provide a directed acyclic graph with $3\epsilon_0$ changed edges). Hence, for the hypothesis $h : \mathbb{R}^d \rightarrow \mathbb{S}_k$, where $h(\mathbf{x})$ is the output of the minimum FAS approximation algorithm with input G (G depends on the input \mathbf{x} , the randomness of the samples and the internal randomness of the $\binom{k}{2}$ calls of the Massart linear classifiers), and the target function $\sigma^*(\mathbf{x})$, we have that

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [\Delta_{KT}(h(\mathbf{x}), \sigma^*(\mathbf{x}))] \leq (\epsilon' + 3\epsilon') = 4\epsilon',$$

which completes the proof, by setting $\epsilon' = \epsilon/4$. \square

Remark 1. *Consider the following variant of the above procedure: compute the $O(k^2)$ linear classifiers with accuracy $\epsilon' = \epsilon/k^2$: If the induced directed graph is acyclic, output the ranking; otherwise, output a random permutation. With probability ϵ , the KT distance will be of order k^2 . Hence, one has to draw in total $O(k^4 d/\epsilon)$ samples to make the expected KT distance roughly $O(\epsilon)$. The algorithm of Theorem 3 improves on this approach.*

A.2 The Proof of Theorem 1: Properly Learning LSFs with Bounded Noise

We first restate the main result of this section.

Theorem 4 (Proper Learning Algorithm). *Fix $\eta \in [0, 1/2)$ and $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2. ProperLSF (Algorithm 2) draws $N = \tilde{O}\left(\frac{d}{\epsilon(1-2\eta)^6} \log(k/\delta)\right)$ samples from \mathcal{D} , runs in $\text{poly}(d, k, 1/\epsilon, \log(1/\delta))$ time and, with probability at least $1 - \delta$, outputs a Linear Sorting function $h : \mathbb{R}^d \rightarrow \mathbb{S}_k$ that is ϵ -close in KT distance to the target.*

We are now ready to provide the proof of our efficient proper learning algorithm for the class of Linear Sorting functions in the presence of bounded noise with respect to the standard Gaussian probability measure.

Proof. As a first step, the algorithm calls the improper learning algorithm ImproperLSF (Algorithm 1) with parameters ϵ, δ and $\eta < 1/2$ and obtains a list of linear classifiers with normal vectors \mathbf{v}_{ij} for $i < j$. The utility of this step implies that, with probability at least $1 - \delta$, each one of the classifiers ϵ -learns the associated true halfspace, i.e., it holds

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\text{sgn}(\mathbf{v}_{ij} \cdot \mathbf{x}) \neq \text{sgn}((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x})] \leq \epsilon,$$

where \mathbf{W}^* is the matrix of the target Linear Sorting function. Without loss of generality, assume that $\|\mathbf{v}_{ij}\|_2 = 1$. In order to make the learner proper, it suffices to solve the following convex program on \mathbf{W} :

$$\text{Find } \mathbf{W} \in \mathbb{R}^{k \times d}, \quad (1)$$

$$\text{such that } (\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{v}_{ij} \geq (1 - \phi) \cdot \|\mathbf{W}_i - \mathbf{W}_j\|_2 \text{ for any } 1 \leq i < j \leq k, \quad (\text{CP}) \quad (2)$$

$$\|\mathbf{W}\|_F \leq 1, \quad (3)$$

for some $\phi \in (0, 1)$ to be decided. The main key ideas are summarized in the next claim.

Claim 5. *The following properties hold true for $\phi = O(\epsilon^2)$ with probability at least $1 - \delta$.*

1. *The convex program 1 is feasible.*
2. *Any solution of the convex program 1 induces an LSF that is ϵ -close in KT distance to the true target $\sigma_{\mathbf{W}^*}(\cdot)$.*
3. *The feasible set of the convex program 1 contains a ball of radius $r = 2^{-\text{poly}(d, k, 1/\epsilon, \log(1/\delta))}$ and is contained in a ball of radius 1. Both balls are with respect to the Frobenius norm.*
4. *The convex program 1 can be solved in time $\text{poly}(d, k, 1/\epsilon, \log(1/\delta))$ using the ellipsoid algorithm.*

Proof of Item 1. First, we can choose the error ϕ so that this convex program is feasible. Let us set $\mathbf{W} = \mathbf{W}^*$, where \mathbf{W}^* is the underlying matrix of the target Linear Sorting function σ^* with $\sigma^*(\mathbf{x}) = \text{argsort}(\mathbf{W}^* \mathbf{x})$. Recall that, by the guarantees of the improper learning algorithm, for the pair $1 \leq i < j \leq k$, it holds

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\text{sgn}(\mathbf{v}_{ij} \cdot \mathbf{x}) \neq \text{sgn}((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x})] \leq \epsilon. \quad (4)$$

Since the standard Gaussian is rotationally symmetric, the angle $\theta(\mathbf{u}, \mathbf{v})$ between two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ is equal to $\pi \cdot \Pr_{\mathbf{x} \sim \mathcal{N}_d} [\text{sgn}(\mathbf{u} \cdot \mathbf{x}) \neq \text{sgn}(\mathbf{v} \cdot \mathbf{x})]$. Hence, using this observation and Equation (4), we get that the angle between the guess vector \mathbf{v}_{ij} and the true normal vector $\mathbf{W}_i^* - \mathbf{W}_j^*$ is

$$\theta(\mathbf{W}_i^* - \mathbf{W}_j^*, \mathbf{v}_{ij}) \leq c \cdot \epsilon,$$

for some constant $c > 0$. For sufficiently small ϵ , this bound implies that the cosine of the above angle is of order $1 - (c\epsilon)^2$ and so the following inequality will hold

$$(\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{v}_{ij} \geq (1 - 2(c\epsilon)^2) \cdot \|\mathbf{W}_i^* - \mathbf{W}_j^*\|_2,$$

since \mathbf{v}_{ij} is unit. Hence, by setting $\phi = 2(c\epsilon)^2$, the convex program with variables $\mathbf{W} \in \mathbb{R}^{k \times d}$ will be feasible; \mathbf{W}^* will be a solution with probability $1 - \delta$, where the randomness is over the output of the algorithm dealing with the Massart linear classifiers. Note that we can assume that $\|\mathbf{W}^*\|_F \leq 1$ without loss of generality, since we can divide each row with the Frobenius norm.

Proof of Item 2. Let $\widetilde{\mathbf{W}}$ be a solution of the convex program. We will make use of the observation that the angle between two vectors is equal to the disagreement of the associated linear threshold functions with respect to the standard normal times π . Observe that any solution $\widetilde{\mathbf{W}}$ to the convex program will satisfy that

$$(\forall i, j) \quad \theta(\mathbf{v}_{ij}, \widetilde{\mathbf{W}}_i - \widetilde{\mathbf{W}}_j) \leq O(\sqrt{\phi}) = c\epsilon.$$

and

$$(\forall i, j) \quad \theta(\mathbf{W}_i^* - \mathbf{W}_j^*, \mathbf{v}_{ij}) \leq \epsilon.$$

This implies that

$$d_{\text{angle}}(\mathbf{W}^*, \widetilde{\mathbf{W}}) \leq c' \epsilon$$

Claim 6. *For the matrices $\mathbf{W}, \mathbf{W}^* \in \mathbb{R}^{k \times d}$, it holds that*

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [\Delta_{\text{KT}}(\sigma_{\mathbf{W}}(\mathbf{x}), \sigma_{\mathbf{W}^*}(\mathbf{x}))] \leq d_{\text{angle}}(\mathbf{W}, \mathbf{W}^*).$$

Proof. We have that

$$\begin{aligned}
\mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [\Delta_{\text{KT}}(\sigma_{\mathbf{W}}(\mathbf{x}), \sigma_{\mathbf{W}^*}(\mathbf{x}))] &= \frac{1}{\binom{k}{2}} \cdot \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} \left[\sum_{1 \leq i < j \leq k} \mathbf{1}\{((\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{x}) ((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x}) < 0\} \right] \\
&= \frac{1}{\binom{k}{2}} \cdot \sum_{1 \leq i < j \leq k} \mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [\text{sgn}(\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{x} \neq \text{sgn}((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x})] \\
&= \frac{1}{\pi} \max_{i,j} \theta(\mathbf{W}_i - \mathbf{W}_j, \mathbf{W}_i^* - \mathbf{W}_j^*) \\
&\leq d_{\text{angle}}(\mathbf{W}, \mathbf{W}^*).
\end{aligned}$$

□

Using the above claim, we get an expected KT distance bound of order $O(\epsilon)$. This gives the desired result.

Proof of Item 3. We will make use of the next lemma.

Lemma 8. Fix $\epsilon, \delta \in (0, 1)$. Let $\mathbf{W}^* \in \mathbb{R}^{k \times d}$ be the true parameter matrix. There exists a matrix $\widetilde{\mathbf{W}}^* \in \mathbb{R}^{k \times d}$ such that, with probability at least $1 - \delta$:

- $\mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [\text{sgn}((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x}) \neq \text{sgn}((\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*) \cdot \mathbf{x})] \leq \epsilon$ for all $i \neq j$, and,
- $\|\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*\|_2 \geq 2^{-\text{poly}(d, k, 1/\epsilon, \log(1/\delta))}$ for any $i \neq j$.

Proof of Lemma 8. The above lemma is a result of the next Appendix A.2.1. In particular, it is a direct implication of Lemma 10 and Corollary 1. □

Note that the above lemma implies that

$$(\forall i, j) \quad \mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [\text{sgn}(\mathbf{v}_{ij} \cdot \mathbf{x}) \neq \text{sgn}((\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*) \cdot \mathbf{x})] \leq 2\epsilon,$$

with probability at least $1 - 2\delta$. Hence, up to constants, the analysis concerning the feasibility of the true matrix \mathbf{W}^* (see Item 1) will still hold for $\widetilde{\mathbf{W}}^*$. From now on we can work with this matrix $\widetilde{\mathbf{W}}^*$ which enjoys the “well-conditionedness” property of the second item of the lemma.

We will use the above lemma in order to prove Item 3 which controls the volume of the feasible region: it states that there exist $0 < r < R$ so that the feasible region of the convex program contains a ball of radius r and is contained in a ball of radius R (where the balls are with respect to the Frobenius norm). Moreover, $r = 2^{-\text{poly}(d, k, 1/\epsilon, \log(1/\delta))}$ and $R = 1$.

For the chosen $\phi \in (0, 1)$, the feasible set contains matrices $\mathbf{W} \in \mathbb{R}^{k \times d}$ that satisfy $\|\mathbf{W} - \widetilde{\mathbf{W}}^*\|_F \leq 2r$, r to be decided. For any $i \neq j$, we have that the following properties hold:

1. $\|\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*\|_2 \geq 2^{-\text{poly}(d, k, 1/\epsilon, \log(1/\delta))}$ (well-conditionedness).
2. $(\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*) \cdot \mathbf{v}_{ij} \geq (1 - \phi) \|\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*\|_2$ (feasibility).
3. $\|\mathbf{W} - \widetilde{\mathbf{W}}^*\|_F \leq 2r$ which implies that $\|\mathbf{W}_i - \widetilde{\mathbf{W}}_i^*\|_2 \leq 2r$ for any $i \in [k]$ (ball around feasible point).
4. $\|\mathbf{v}_{ij}\|_2 = 1$.

Our goal is to prove that for a matrix in the above ball it holds $(\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{v}_{ij} \geq (1 - \phi) \|\mathbf{W}_i - \mathbf{W}_j\|_2$.

We have that

$$\begin{aligned}
(\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*) \cdot \mathbf{v}_{ij} &= (\widetilde{\mathbf{W}}_i^* - \mathbf{W}_i) \cdot \mathbf{v}_{ij} + (\mathbf{W}_j - \widetilde{\mathbf{W}}_j^*) \cdot \mathbf{v}_{ij} + (\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{v}_{ij} \\
&\leq \|\widetilde{\mathbf{W}}_i^* - \mathbf{W}_i\|_2 + \|\mathbf{W}_j - \widetilde{\mathbf{W}}_j^*\|_2 + (\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{v}_{ij} \\
&\leq 4r + (\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{v}_{ij}.
\end{aligned}$$

More to that

$$\begin{aligned}
\|\mathbf{W}_i - \mathbf{W}_j\|_2 &= \|\mathbf{W}_i - \widetilde{\mathbf{W}}_i^* + \widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^* + \widetilde{\mathbf{W}}_j^* - \mathbf{W}_j\|_2 \\
&\leq \|\mathbf{W}_i - \widetilde{\mathbf{W}}_i^*\|_2 + \|\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*\|_2 + \|\widetilde{\mathbf{W}}_j^* - \mathbf{W}_j\|_2 \\
&\leq 4r + \|\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*\|_2,
\end{aligned}$$

and similarly: $\|\mathbf{W}_i - \mathbf{W}_j\|_2 \geq \|\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*\|_2 - 4r$.

Combining the above inequalities, we get that

$$\begin{aligned}
(\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{v}_{ij} &\geq (\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*) \cdot \mathbf{v}_{ij} - 4r \\
&\geq (1 - \phi) \|\widetilde{\mathbf{W}}_i^* - \widetilde{\mathbf{W}}_j^*\|_2 - 4r \\
&\geq (1 - \phi) (\|\mathbf{W}_i - \mathbf{W}_j\|_2 - 4r) - 4r \\
&= (1 - \phi) \|\mathbf{W}_i - \mathbf{W}_j\|_2 - 8r.
\end{aligned}$$

We pick r sufficiently small and of order $2^{-\text{poly}(d, k, 1/\epsilon, \log(1/\delta))}$ and get that \mathbf{W} is a feasible solution of the convex program. Moreover, we can select $R = 1$ since $\|\widetilde{\mathbf{W}}^*\|_F = 1$ without loss of generality, since we can normalize the row differences of $\widetilde{\mathbf{W}}^*$ with the norm $\|\widetilde{\mathbf{W}}^*\|_F$.

Proof of Item 4. We apply the ellipsoid algorithm in order to solve the convex program 1 and compute a matrix $\widetilde{\mathbf{W}} \in \mathbb{R}^{k \times d}$. The algorithm ProperLSF outputs the linear sorting function $h(\cdot) = \sigma_{\widetilde{\mathbf{W}}}(\cdot)$.

Lemma 9 (Efficiency of the Ellipsoid Algorithm [Vis21]). *Suppose that $P \subseteq \mathbb{R}^d$ is a full-dimensional polytope that is contained in a d -dimensional Euclidean ball of radius $R > 0$ and contains a d -dimensional Euclidean ball of radius $r > 0$. Then, the ellipsoid method outputs a point $\tilde{\mathbf{x}} \in P$ after $O(d^2 \log(R/r))$ iterations. Moreover, every iteration can be implemented in $O(d^2 + T_{\text{sep}})$ time, where T_{sep} is the time required to answer a single query by the separation oracle.*

Assume that Item 3 holds true. Then the algorithm can be used with $r = 2^{-\text{poly}(d, k, 1/\epsilon, \log(1/\delta))}$ and $R = 1$. Hence, the ellipsoid algorithm will provide in time $\text{poly}(d, k, 1/\epsilon, \log(1/\delta))$ a point $\widetilde{\mathbf{W}}$ that lies in the feasible region of the convex program 1⁵. □

Remark 2. *We remark that both the improper (Algorithm 1) and the proper (Algorithm 2) learning algorithms hold for the more general case where the \mathbf{x} -marginal lies in the class of isotropic log-concave distributions [LV07]: A distribution \mathcal{D}_x lies inside the class of isotropic log-concave distributions \mathcal{F}_{LC} over \mathbb{R}^d if \mathcal{D}_x has a probability density function f over \mathbb{R}^d such that $\log f$ is concave, its mean is zero, and its covariance is identity, i.e., $\mathbf{E}_{\mathbf{x} \sim \mathcal{D}_x}[\mathbf{x}\mathbf{x}^\top] = \mathbf{I}$.*

A.2.1 The proof of Lemma 8

We provide the following result.

Lemma 10. *Fix $\epsilon, \delta \in (0, 1)$. Let $\mathbf{W}^* \in \mathbb{R}^{k \times d}$ be the true parameter matrix. There exists a matrix $\mathbf{W} \in \mathbb{R}^{k \times d}$ such that, with probability at least $1 - \delta$:*

- $\Pr_{\mathbf{x} \sim \mathcal{N}_d}[\text{sgn}((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x}) \neq \text{sgn}((\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{x})] \leq \epsilon$ for all $i \neq j$, and,
- The bit complexity of \mathbf{W} is $\text{poly}(k, d, 1/\epsilon, \log(1/\delta))$

Proof. The matrix \mathbf{W} will be the output of a linear program that can be used to learn the LSF $\sigma_{\mathbf{W}^*}(\cdot)$ in the noiseless setting.

⁵We remark that the runtime will also depend on the time required to answer a single query by the separation oracle. We assume that this time is polynomial in the parameters of our problem and we opt not to track these details in this work.

Consider the unit sphere \mathcal{S}^{d-1} and a δ_0 -cover of the unit sphere with parameter $\delta_0 > 0$ to be decided. For any sample $(\mathbf{x}, \pi) \sim \mathcal{D}$ of the 0-noisy linear label ranking distribution, i.e., $\mathbf{x} \sim \mathcal{N}_d$ and $\pi = \sigma_{\mathbf{W}^*}(\mathbf{x})$, we consider the rounded sample $(\tilde{\mathbf{x}}, \pi)$ where $\tilde{\mathbf{x}}$ is obtained by first projecting $\mathbf{x} \in \mathbb{R}^d$ to \mathcal{S}^{d-1} and then by obtaining the closest point of $\hat{\mathbf{x}}$ in the cover. The cover's size is $O(1/\delta_0)^d$.

Let us fix $1 \leq i < j \leq k$ and set $y_{ij} = \text{sgn}(\pi(i) - \pi(j))$. For a training set $\{(\mathbf{x}^{(t)}, \pi^{(t)})\}_{t \in [N]}$ of size N , we create the following linear system L_{ij} with variables $\mathbf{W} \in \mathbb{R}^{k \times d}$:

$$y_{ij}^{(t)} (\mathbf{W}_i - \mathbf{W}_j) \cdot \tilde{\mathbf{x}}^{(t)} \geq 0, t \in [N] \quad (L_{ij}).$$

Consider the concatenation of the linear systems $L = \cup_{i < j} L_{ij}$. The number of equations in the linear system of equations L is $N \cdot \binom{k}{2}$.

We first have to show that, with high probability, the system L is feasible, i.e., there exists \mathbf{W} that satisfies the system's equations. Note that if we replace $\tilde{\mathbf{x}}^{(t)}$ with the original points $\mathbf{x}^{(t)}$, the true matrix \mathbf{W}^* is a solution to the system. We now have to study the rounded linear system.

Claim 7. *The (rounded) linear system L is feasible with high probability.*

Proof. In order to show the feasibility of L , we will use the anti-concentration properties of the Gaussian.

Fact 1 ([DKM05]). *Let \mathcal{P} be the standard normal distribution over \mathbb{R}^d . For any fixed unit vector $\mathbf{a} \in \mathbb{R}^d$ and any $\gamma \leq 1$,*

$$\gamma/4 \leq \Pr_{\mathbf{x} \sim \mathcal{P}} \left[|\mathbf{a} \cdot \frac{\mathbf{x}}{\|\mathbf{x}\|_2}| \leq \frac{\gamma}{\sqrt{d}} \right] \leq \gamma.$$

Let us focus on the pair $1 \leq i < j \leq k$. We first observe that scaling all samples to lie on the unit sphere does not affect the feasibility of the system. It suffices to focus on that single halfspace with normal vector $\mathbf{v}_{ij} = \mathbf{W}_i^* - \mathbf{W}_j^* \in \mathbb{R}^d$ and consider the probability of the event that the collection of the N rounded points $\{\tilde{\mathbf{x}}^{(t)}\}_t$ with labels $\{y_{ij}^{(t)}\}_t$, that come from N Gaussian vectors $\{\mathbf{x}^{(t)}\}_t$ which are linearly separable (with labels $\{y_{ij}^{(t)}\}_t$), becomes non-linearly separable. For this it suffices to control the probability that the rounding procedure flips the label of the data point. Using the union bound, we have that, if the rounding has accuracy δ_0 , the described bad event has probability

$$\Pr_{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)} \sim \mathcal{N}_d} [\exists t \in [N] : \text{sgn}(\mathbf{v}_{ij} \cdot \tilde{\mathbf{x}}^{(t)}) \neq \text{sgn}(\mathbf{v}_{ij} \cdot \mathbf{x}^{(t)})] \leq N \cdot \Pr_{\mathbf{x} \sim \mathcal{N}_d} [|\mathbf{v}_{ij} \cdot \mathbf{x} / \|\mathbf{x}\|_2| \leq 2\delta_0] \leq N \cdot O(\delta_0 \sqrt{d}),$$

where we remark that the first event is scale invariant and so we can assume that the normal vector is unit, the first inequality follows from the fact that it suffices to control the mass assigned to a strip of width $2\delta_0$ (due to the discretization) and the second inequality follows from Fact 1. We now have to select the discretization. Let $\delta \in (0, 1)$. By choosing $\delta_0 = O(\frac{\delta}{N\sqrt{dk^2}})$, the bad event for all the pairs $i < j$ occurs with probability at most δ , i.e., with probability at least $1 - \delta$, each one of the N drawn i.i.d. samples does not fall in any one of the $\binom{k}{2}$ “bad” strips. \square

We can now consider the case that the system L is feasible (with the target matrix \mathbf{W}^* being a feasible point) that occurs with probability $1 - \delta$. The class of homogenous halfspaces in d dimensions has VC dimension d ; therefore, the sample complexity of learning halfspaces using ERM is $O((d + \log(1/\delta))/\epsilon)$. Moreover, in the realizable case, we can implement the ERM using e.g., linear programming and find a solution in $\text{poly}(d, 1/\epsilon, \log(1/\delta))$ time. We next focus on the quality of the solution which will give the desired sample complexity.

Claim 8. *Assume that the algorithm draws $N = \tilde{O}(\frac{d + \log(k/\delta)}{\epsilon})$ i.i.d. samples of the form (\mathbf{x}, π) with $\mathbf{x} \sim \mathcal{N}_d$ and $\pi = \sigma_{\mathbf{W}^*}(\mathbf{x})$. For any $i \neq j$ and with probability at least $1 - 2\delta$, the solution \mathbf{W} of the linear system L satisfies*

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\text{sgn}((\mathbf{W}_i^* - \mathbf{W}_j^*) \cdot \mathbf{x}) \neq \text{sgn}((\mathbf{W}_i - \mathbf{W}_j) \cdot \mathbf{x})] \leq \epsilon.$$

Proof. Since the matrix \mathbf{W} satisfies the sub-system L_{ij} , the result follows using a union bound on the events that (i) the linear system is feasible and (ii) the ERM is a successful PAC learner. \square

Claim 9. Consider the solution \mathbf{W} of the linear system. Then, \mathbf{W} has bounded bit complexity of order $\text{poly}(d, k, 1/\epsilon, \log(1/\delta))$.

Proof. We will make use of the following result that relates the size of the input and the output of a linear program using Cramer's rule.

Lemma 11 ([Sch98, Pap81]). Let $\mathbf{A} \in \mathbb{Z}^{m \times n}$, $\mathbf{b} \in \mathbb{Z}^m$, $\mathbf{c} \in \mathbb{Z}^n$. Consider a linear program $\min \mathbf{c} \cdot \mathbf{x}$ subject to $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ and $\mathbf{x} \geq \mathbf{0}$. Let U be the maximum size of A_{ij} , b_i , c_j . The output of the linear program has size $O(m(nU + n \log(n)))$ bits.

We will apply the above lemma (which holds even by dropping the constraint $\mathbf{x} \geq \mathbf{0}$) to our setting where $\mathbf{A}\mathbf{w} \geq \mathbf{0}$ where $\mathbf{w} = (\mathbf{W}_i)_{i \in [k]} \in \mathbb{Q}^{kd}$, i.e., \mathbf{w} is the vectorization of the matrix \mathbf{W} . Moreover, \mathbf{A} is the matrix containing the N (rounded) Gaussian samples $\tilde{\mathbf{x}}^{(t)}$. We have that the matrix \mathbf{A} has dimension $N \binom{k}{2} \times kd$ and each entry A_{ij} is an integer and has size at most $U = \text{poly}(d, k)$ (since the samples are rounded on the δ_0 -cover of the sphere. Recall that the labels $y_{ij}^{(t)} \in \{-1, +1\}$ and $\tilde{\mathbf{x}}^{(t)}$ lie in the unit sphere. In particular, each row of the matrix \mathbf{A} has $2d$ non-zero entries and is associated with a tuple (i, j, t) for $1 \leq i < j \leq k$ and $t \in [N]$. Then, it holds that the output has size at most $O(Nk^2(dU + dk \log(dk)))$ bits. So, we get that the output \mathbf{W} can be described using at most $\text{poly}(d, k, 1/\epsilon, U, \log(1/\delta)) = \text{poly}(d, k, 1/\epsilon, \log(1/\delta))$ bits (due to the size of the entries of the matrix \mathbf{A}). \square

Combining the above claims, we conclude the proof. \square

As a corollary of the bounded bit complexity, we obtain the following key result.

Corollary 1. Let $\epsilon > 0$. Assume that $\mathbf{W} \in \mathbb{R}^{k \times d}$ has bit complexity at most $\text{poly}(d, k, 1/\epsilon, \log(1/\delta))$. Then, for any $i, j \in [k]$ with $i \neq j$, it holds that $\|\mathbf{W}_i - \mathbf{W}_j\|_2 > 2^{-\text{poly}(d, k, 1/\epsilon, \log(1/\delta))}$.

Proof. First, we can assume that $\mathbf{W}_i \neq \mathbf{W}_j$ for any $i \neq j$; in case of equal rows, we obtain a low-dimensional instance. Then, since any vector \mathbf{W}_i has bounded bit complexity, we have that the difference of any two such vectors, provided that it is non-zero, has a lower bound in its norm, i.e., $\|\mathbf{W}_i - \mathbf{W}_j\|_2 > 2^{-\text{poly}(d, k, 1/\epsilon, \log(1/\delta))}$ for any $i, j \in [k]$. \square

B Learning in Top-1 Disagreement from Label Rankings

Let us set $\sigma_1(\mathbf{W}\mathbf{x}) = \text{argmax}_{i \in [k]} \mathbf{W}_i \cdot \mathbf{x}$ for $\mathbf{x} \in \mathbb{R}^d$. The main result of this section follows.

Theorem 5 (Proper Top-1 Learning Algorithm). Fix $\eta \in [0, 1/2)$ and $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2. There exists an algorithm that draws $N = O\left(\frac{dk\sqrt{\log k}}{\epsilon(1-2\eta)^6} \log(k/\delta)\right)$ samples from \mathcal{D} , runs in $\text{poly}(N)$ time and, with probability at least $1 - \delta$, outputs a Linear Sorting function $h : \mathbb{R}^d \rightarrow \mathbb{S}_k$ that is ϵ -close in top-1 disagreement to the target.

Proof. Note that the `MassartLTF` algorithm (see Lemma 6) has the guarantee that it returns a vector \mathbf{w} so that

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\text{sgn}(\mathbf{w} \cdot \mathbf{x}) \neq \text{sgn}(\mathbf{w}^* \cdot \mathbf{x})] \leq \epsilon,$$

with probability $1 - \delta$, where \mathbf{w}^* is the target normal vector. Since the above misclassification probability with respect to \mathcal{N}_d is directly connected with the angle $\theta(\mathbf{w}, \mathbf{w}^*)$, we get that we can control the angle between \mathbf{w} and \mathbf{w}^* efficiently. Moreover, in our setting, for a matrix $\mathbf{W} \in \mathbb{R}^{k \times d}$, there exist $\binom{k}{2}$ homogeneous halfspaces with normal vectors $\mathbf{W}_i - \mathbf{W}_j$ and so we can control the angles $\theta(\mathbf{W}_i - \mathbf{W}_j, \mathbf{W}_i^* - \mathbf{W}_j^*)$. In order to deduce the sample complexity bound of Theorem 5, we show the next lemma which essentially bounds the top-1 misclassification error using the angles of these $O(k^2)$ halfspaces. We apply Lemma 12 with $\mathbf{U} = \mathbf{W}$ and $\mathbf{V} = \mathbf{W}^*$ and so we can take $\epsilon' = \epsilon/(k\sqrt{\log k})$ and invoke the proper learning algorithm of Algorithm 2. This completes the proof. \square

We continue with the proof of our key lemma.

Lemma 12 (Misclassification Error). *Consider two matrices $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{k \times d}$ and let \mathcal{N}_d be the standard Gaussian in d dimensions. We have that*

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\sigma_1(\mathbf{U}\mathbf{x}) \neq \sigma_1(\mathbf{V}\mathbf{x})] \leq c \cdot k \cdot \sqrt{\log k} \cdot \max_{i \neq j} \theta(\mathbf{U}_i - \mathbf{U}_j, \mathbf{V}_i - \mathbf{V}_j),$$

where $c > 0$ is some universal constant.

Proof. We have that

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\sigma_1(\mathbf{U}\mathbf{x}) \neq \sigma_1(\mathbf{V}\mathbf{x})] = \sum_{i \in [k]} \Pr_{\mathbf{x} \sim \mathcal{N}_d} [\sigma_1(\mathbf{U}\mathbf{x}) = i, \sigma_1(\mathbf{V}\mathbf{x}) \neq i].$$

We have that $\mathcal{C}_{\mathbf{U}}^{(i)} = \mathbf{1}\{\mathbf{x} : \sigma_1(\mathbf{U}\mathbf{x}) = i\} = \prod_{j \neq i} \mathbf{1}\{(\mathbf{U}_i - \mathbf{U}_j) \cdot \mathbf{x} \geq 0\}$ is the set indicator of a homogeneous polyhedral cone as the intersection of $k - 1$ homogeneous halfspaces. Similarly, we consider the cone $\mathcal{C}_{\mathbf{V}}^{(i)} = \{\mathbf{x} : \sigma_1(\mathbf{V}\mathbf{x}) = i\}$. Hence, we have that $\{\mathbf{x} : \sigma_1(\mathbf{V}\mathbf{x}) \neq i\}$ is the complement of a homogeneous polyhedral cone. Let us define $C_{\mathbf{U}}^{(i)} : \mathbb{R}^d \mapsto \{0, 1\}$ and $C_{\mathbf{V}}^{(i)} : \mathbb{R}^d \mapsto \{0, 1\}$ be the associated indicator functions of the two cones. We have that

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\sigma_1(\mathbf{U}\mathbf{x}) = i, \sigma_1(\mathbf{V}\mathbf{x}) \neq i] = \Pr_{\mathbf{x} \sim \mathcal{N}_d} [C_{\mathbf{U}}^{(i)}(\mathbf{x}) = 1, C_{\mathbf{V}}^{(i)}(\mathbf{x}) = 0].$$

Finally, we have that

$$C_{\mathbf{U}}^{(i)} \cap \left(C_{\mathbf{V}}^{(i)}\right)^c = C_{\mathbf{U}}^{(i)} \setminus C_{\mathbf{V}}^{(i)} \subseteq C_{\mathbf{U}}^{(i)} \setminus C_{\mathbf{V}}^{(i)} \cup C_{\mathbf{V}}^{(i)} \setminus C_{\mathbf{U}}^{(i)}.$$

We can hence apply Lemma 13 for the cones $C_{\mathbf{U}}^{(i)}, C_{\mathbf{V}}^{(i)}$ for each $i \in [k]$. \square

Lemma 13 (Cone Disagreement). *Let $C_1 : \mathbb{R}^d \mapsto \{0, 1\}$ be the indicator function of the homogeneous polyhedral cone defined by the k unit vectors $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^d$, i.e., $C_1(\mathbf{x}) = \prod_{i=1}^k \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \geq 0\}$. Similarly, define $C_2 : \mathbb{R}^d \mapsto \{0, 1\}$ to be the homogeneous polyhedral cone with normal vectors $\mathbf{u}_1, \dots, \mathbf{u}_k$. It holds that*

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [C_1(\mathbf{x}) \neq C_2(\mathbf{x})] \leq c \sqrt{\log(k)} \max_{i \in [k]} \theta(\mathbf{v}_i, \mathbf{u}_i),$$

where $c > 0$ is some universal constant.

Proof. To simplify notation, denote $\theta = \max_{i \in [k]} \theta(\mathbf{v}_i, \mathbf{u}_i)$. We first observe that it suffices to prove the upper bound on the probability of $C_1(\mathbf{x}) \neq C_2(\mathbf{x})$ for sufficiently small values of θ . Indeed, if we have that the bound is true for θ smaller than some θ_0 we can then form a path of sufficiently large length N (in particular we need $\theta/N \leq \theta_0$) starting from the vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ to the final vectors $\mathbf{u}_1, \dots, \mathbf{u}_k$, where at each step we only rotate the vectors by at most $\theta/N \leq \theta_0$. By the triangle inequality, we immediately obtain that the probability that $C_1(\mathbf{x}) \neq C_2(\mathbf{x})$ is at most equal to the sum of the probabilities of the intermediate steps which is at most $\sum_{i=1}^N c \sqrt{\log(k)} \frac{\theta}{N} = c \sqrt{\log(k)} \theta$. Notice in the above argument the constant θ_0 can be arbitrarily small and may also depend on k and d .

We define the indicator of the positive orthant in k dimensions to be $R(\mathbf{t}) = \prod_{i=1}^k \mathbf{1}\{t_i \geq 0\}$. Using this notation, we have that the cone indicator can be written as $C_1(\mathbf{x}) = R(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) = R(\mathbf{V}\mathbf{x})$, where \mathbf{V} is the $k \times d$ matrix whose i -th row is the vector \mathbf{v}_i . Moreover, we define the i -th face of the cone $R(\mathbf{V}\mathbf{x})$ to be

$$F_i(\mathbf{V}\mathbf{x}) = R(\mathbf{V}\mathbf{x}) \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} = 0\}.$$

We will first handle the case where only one of the normal vectors \mathbf{v}_i changes. We show the following claim.

Claim 10. *Let $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^d$ and $\mathbf{r} \in \mathbb{R}^d$ with $\theta(\mathbf{v}_1, \mathbf{r}) \leq \theta$ for some sufficiently small $\theta \in (0, \pi/2)$. It holds that*

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \neq R(\mathbf{r} \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})] \leq c \cdot \theta \cdot \Gamma(F_1) \sqrt{\log \left(\frac{1}{\Gamma(F_1)} + 1 \right)},$$

where F_1 is the face with $\mathbf{v}_1 \cdot \mathbf{x} = 0$ of the cone $R(\mathbf{V}\mathbf{x})$ and c is some universal constant.

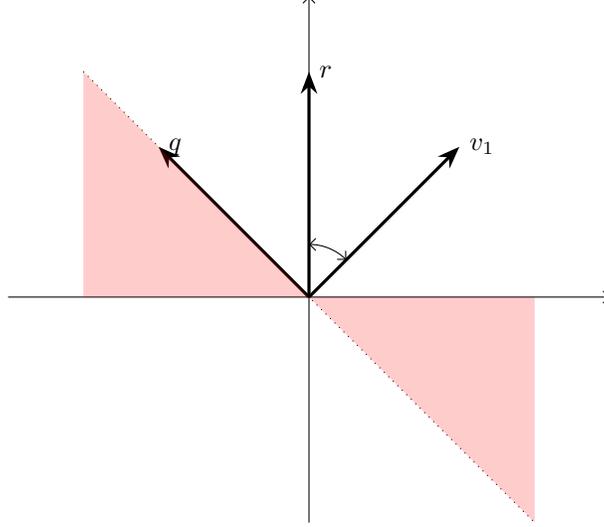


Figure 1: The vectors \mathbf{r} , \mathbf{v}_1 and \mathbf{q} and the disagreement region of the halfspaces with normal vectors \mathbf{r} and \mathbf{v}_1 .

Proof. We have

$$\begin{aligned} & \Pr_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \neq R(\mathbf{r} \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})] \\ &= \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [|R(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) - R(\mathbf{r} \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})|] \\ &= \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) |\mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \geq 0\} - \mathbf{1}\{\mathbf{r} \cdot \mathbf{x} \geq 0\}|] . \end{aligned}$$

We have that $|\mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \geq 0\} - \mathbf{1}\{\mathbf{r} \cdot \mathbf{x} \geq 0\}| = \mathbf{1}\{(\mathbf{v}_1 \cdot \mathbf{x})(\mathbf{r} \cdot \mathbf{x}) < 0\}$, i.e., this is the event that the halfspaces $\mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \geq 0\}$ and $\mathbf{1}\{\mathbf{r} \cdot \mathbf{x} \geq 0\}$ disagree. Let \mathbf{q} be the normalized projection of \mathbf{r} onto the orthogonal complement of \mathbf{v}_1 , i.e., $\mathbf{q} = \text{proj}_{\mathbf{v}_1^\perp} \mathbf{r} / \|\text{proj}_{\mathbf{v}_1^\perp} \mathbf{r}\|_2$. We have that \mathbf{v}_1 and \mathbf{q} is an orthonormal basis of the subspace spanned by the vectors \mathbf{v}_1 and \mathbf{r} . We have that $\mathbf{r} = \cos \theta(\mathbf{v}_1, \mathbf{r})\mathbf{v}_1 + \sin \theta(\mathbf{v}_1, \mathbf{r})\mathbf{q}$. Moreover, we have that the region $(\mathbf{v}_1 \cdot \mathbf{x})(\mathbf{r} \cdot \mathbf{x}) < 0$ is equal to

$$\{0 < \mathbf{v}_1 \cdot \mathbf{x} < -(\mathbf{q} \cdot \mathbf{x}) \tan \theta(\mathbf{v}_1, \mathbf{r})\} \cup \{-(\mathbf{q} \cdot \mathbf{x}) \tan \theta(\mathbf{v}_1, \mathbf{r}) < \mathbf{v}_1 \cdot \mathbf{x} < 0\} .$$

Thus, we have that the disagreement region $(\mathbf{v}_1 \cdot \mathbf{x})(\mathbf{r} \cdot \mathbf{x}) < 0$ is a subset of the region $\{|\mathbf{v}_1 \cdot \mathbf{x}| \leq |\mathbf{q} \cdot \mathbf{x}| \tan \theta(\mathbf{v}_1, \mathbf{r})\}$. Since $\tan \theta(\mathbf{v}_1, \mathbf{r}) \leq \theta$ and we have that θ is sufficiently small we can also replace the above region by the larger region: $\{|\mathbf{v}_1 \cdot \mathbf{x}| \leq 2\theta|\mathbf{q} \cdot \mathbf{x}|\}$. Therefore, we have

$$\begin{aligned} & \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \mathbf{1}\{(\mathbf{v}_1 \cdot \mathbf{x})(\mathbf{r} \cdot \mathbf{x}) < 0\}] \\ & \leq \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \mathbf{1}\{|\mathbf{v}_1 \cdot \mathbf{x}| \leq 2\theta|\mathbf{q} \cdot \mathbf{x}|\}] . \end{aligned}$$

The derivative of the above expression with respect to θ is equal to

$$\mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} \left[R(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \delta \left(\frac{|\mathbf{v}_1 \cdot \mathbf{x}|}{2|\mathbf{q} \cdot \mathbf{x}|} - \theta \right) \right] ,$$

where $\delta(t)$ is the Dirac delta function. At $\theta = 0$ and using the property that $\delta(t/a) = a\delta(t)$, we have that the above derivative is equal to

$$2 \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| \delta(|\mathbf{v}_1 \cdot \mathbf{x}|)] .$$

Notice that, if we did not have the term $|\mathbf{q} \cdot \mathbf{x}|$, the above expression would be exactly equal to two times the Gaussian surface area of the face with $\mathbf{v}_1 \cdot \mathbf{x} = 0$, i.e., it would be equal to $2\Gamma(F_1)$. We now show that this extra term of $|\mathbf{q} \cdot \mathbf{x}|$ can only increase the above surface integral by at most a

logarithmic factor. We have that

$$\begin{aligned}
\mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) | \mathbf{q} \cdot \mathbf{x} | \delta(|\mathbf{v}_1 \cdot \mathbf{x}|)] &= \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| d\mu(\mathbf{x}) \\
&\leq \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| \mathbf{1}\{|\mathbf{q} \cdot \mathbf{x}| \leq \xi\} d\mu(\mathbf{x}) + \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| \mathbf{1}\{|\mathbf{q} \cdot \mathbf{x}| \geq \xi\} d\mu(\mathbf{x}) \\
&\leq \xi \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) d\mu(\mathbf{x}) + \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| \mathbf{1}\{|\mathbf{q} \cdot \mathbf{x}| \geq \xi\} d\mu(\mathbf{x}),
\end{aligned}$$

where $d\mu(\mathbf{x})$ is the standard surface measure in \mathbb{R}^d . The first term above is exactly equal to the Gaussian surface area of the face F_1 . To bound from above the second term we can use the fact that the face F_1 is a subset of the hyperplane $\mathbf{v}_1 \cdot \mathbf{x} = 0$, i.e., it holds that $F_1 \subseteq \{\mathbf{x} : |\mathbf{v}_1 \cdot \mathbf{x}| = 0\}$. To simplify notation we may assume that $\mathbf{v}_1 = \mathbf{e}_1$ and $\mathbf{q} = \mathbf{e}_2$ (recall that \mathbf{v}_1 and \mathbf{q} are orthogonal unit vectors), and in this case we obtain

$$\begin{aligned}
\int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| \mathbf{1}\{|\mathbf{q} \cdot \mathbf{x}| \geq \xi\} d\mu(\mathbf{x}) &\leq \int_{x_1=0} \phi_d(\mathbf{x}) |x_2| \mathbf{1}\{|x_2| \geq \xi\} d\mu(\mathbf{x}) \\
&= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} |x_2| \mathbf{1}\{|x_2| \geq \xi\} \frac{e^{-x_2^2/2}}{\sqrt{2\pi}} dx_2 \\
&= \frac{1}{\pi} e^{-\xi^2/2}.
\end{aligned}$$

Combining the above bounds we obtain that the derivative with respect to θ of the expression $\mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \mathbf{1}\{|\mathbf{v}_1 \cdot \mathbf{x}| \leq 2\theta |\mathbf{q} \cdot \mathbf{x}|\}]$ is equal to

$$\frac{d}{d\theta} \left(\mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \mathbf{1}\{|\mathbf{v}_1 \cdot \mathbf{x}| \leq 2\theta |\mathbf{q} \cdot \mathbf{x}|\}] \right) \Big|_{\theta=0} \leq 2\xi \Gamma(F_1) + \frac{2e^{-\xi^2/2}}{\pi}.$$

By picking $\xi = \sqrt{2 \log(1 + 1/\Gamma(F_1))}$, the result follows since up to introducing $o(\theta)$ error we can bound the term $\mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \neq R(\mathbf{r} \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})]$ by its derivative with respect to θ (evaluated at 0) times θ . \square

We can complete the proof of Lemma 13 using Claim 10. In order to bound the disagreement of the cones C_1 and C_2 we can start from C_1 and change one of its vectors at a time so that we can use Claim 10 that can handle this case. For example, at the first step, we can swap \mathbf{v}_1 for \mathbf{u}_1 and use the triangle inequality to obtain that

$$\begin{aligned}
\mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [C_1(\mathbf{x}) \neq C_2(\mathbf{x})] &\leq \mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \neq R(\mathbf{u}_1 \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})] \\
&\quad + \mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{u}_1 \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \neq R(\mathbf{u}_1 \cdot \mathbf{x}, \mathbf{u}_2 \cdot \mathbf{x}, \dots, \mathbf{u}_k \cdot \mathbf{x})] \\
&\leq c \cdot \theta \Gamma(F_1) \sqrt{\log(1/\Gamma(F_1) + 1)} \\
&\quad + \mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{u}_1 \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \neq R(\mathbf{u}_1 \cdot \mathbf{x}, \mathbf{u}_2 \cdot \mathbf{x}, \dots, \mathbf{u}_k \cdot \mathbf{x})],
\end{aligned}$$

where $F_1 = F_1(\mathbf{V}\mathbf{x})$ is the face with $\mathbf{v}_1 \cdot \mathbf{x} = 0$ of the cone C_1 . Notice that we have replaced \mathbf{v}_1 by \mathbf{u}_1 in the above bound. Our plan is to use the triangle inequality and continue replacing the vectors of C_1 by the vectors of C_2 sequentially. To make this formal we define the matrix $\mathbf{A}^{(i)} \in \mathbb{R}^{k \times d}$ whose first $i-1$ rows are the vectors $\mathbf{u}_1, \dots, \mathbf{u}_{i-1}$ and its last $k-i+1$ rows are the vectors $\mathbf{v}_i, \dots, \mathbf{v}_k$, i.e.,

$$\mathbf{A}_j^{(i)} = \begin{cases} \mathbf{u}_j & \text{if } 1 \leq j \leq i-1, \\ \mathbf{v}_j & \text{if } i \leq j \leq k. \end{cases}$$

Notice that $\mathbf{A}^{(1)} = \mathbf{V}$ and $\mathbf{A}^{(k+1)} = \mathbf{U}$. Using the triangle inequality we obtain that

$$\mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [C_1(\mathbf{x}) \neq C_2(\mathbf{x})] \leq \sum_{i=1}^k \mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [R(\mathbf{A}^{(i)}\mathbf{x}) \neq R(\mathbf{A}^{(i+1)}\mathbf{x})].$$

Since the matrices $\mathbf{A}^{(i)}$ and $\mathbf{A}^{(i+1)}$ only differ on one row, we can use Claim 10 to obtain the following bound:

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [C_1(\mathbf{x}) \neq C_2(\mathbf{x})] \leq c \cdot \theta \cdot \sum_{i=1}^k \Gamma(F_i(\mathbf{A}^{(i)} \mathbf{x})) \sqrt{\frac{1}{\Gamma(F_i(\mathbf{A}^{(i)} \mathbf{x}))} + 1}.$$

We now observe that the Gaussian surface area $\Gamma(F_i(\mathbf{A}^{(i)} \mathbf{x}))$ is a continuous function of the matrix $\mathbf{A}^{(i)}$. By flattening the matrix $\mathbf{A}^{(i)}$ (since it is isomorphic to a vector $\mathbf{z} \in \mathbb{R}^{n^2}$) and letting $S_{\mathbf{z}}$ be the induced surface $\{\mathbf{x} : R(\mathbf{A}^{(i)} \mathbf{x}) = 1 \wedge \mathbf{v}_i \cdot \mathbf{x} = 0\}$, it suffices to show that

$$\lim_{\mathbf{w} \rightarrow \mathbf{z}} \int \phi_n(\mathbf{x}) \mathbf{1}\{\mathbf{x} \in S_{\mathbf{w}}\} d\mu(\mathbf{x}) = \int \phi_n(\mathbf{x}) \mathbf{1}\{\mathbf{x} \in S_{\mathbf{z}}\} d\mu(\mathbf{x}),$$

by the smoothness of the surface $S_{\mathbf{z}}$. Consider a sequence of functions (g_m) and vectors (\mathbf{w}_m) so that $g_m(\mathbf{x}) = \phi_n(\mathbf{x}) \mathbf{1}\{\mathbf{x} \in S_{\mathbf{w}_m}\}$ and $\lim_{m \rightarrow \infty} \mathbf{w}_m = \mathbf{z}$. Note that $|g_m(\mathbf{x})| \leq 1$ everywhere. Hence, by the dominated convergence theorem, we have that

$$\lim_{m \rightarrow \infty} \int g_m(\mathbf{x}) d\mu(\mathbf{x}) = \int \lim_{m \rightarrow \infty} g_m(\mathbf{x}) d\mu(\mathbf{x}) = \int \phi_n(\mathbf{x}) \lim_{m \rightarrow \infty} \mathbf{1}\{\mathbf{x} \in S_{\mathbf{w}_m}\} d\mu(\mathbf{x}).$$

Since the sequence consists of smooth surfaces, we have that $\lim_{m \rightarrow \infty} \mathbf{1}\{\mathbf{x} \in S_{\mathbf{w}_m}\} = \mathbf{1}\{\mathbf{x} \in S_{\mathbf{z}}\}$ and so the Gaussian surface area is continuous with respect to the matrix $\mathbf{A}^{(i)}$ for any $i \in [k]$.

Also, as $\theta \rightarrow 0$, we have that $\mathbf{A}^{(i)} \rightarrow \mathbf{V}$. This is because the sequence of matrices $\mathbf{A}^{(i)}$ depends only on the vectors \mathbf{u}_j and \mathbf{v}_j for $j \in [k]$ and the following two properties hold true: $\theta = \max_{j \in [k]} \theta(\mathbf{v}_j, \mathbf{u}_j)$ and all the vectors are unit. Hence, as θ tends to zero, they tend to become the same vectors and so any matrix $\mathbf{A}^{(i)}$ tends to become \mathbf{V} . Therefore, taking this limit we obtain that for $\theta \rightarrow 0$ it holds that

$$\lim_{\theta \rightarrow 0} \frac{\Pr_{\mathbf{x} \sim \mathcal{N}_d} [C_1(\mathbf{x}) \neq C_2(\mathbf{x})]}{\theta} \leq c \cdot \sum_{i=1}^k \Gamma(F_i(\mathbf{V} \mathbf{x})) \sqrt{\log(1/\Gamma(F_i(\mathbf{V} \mathbf{x}))) + 1}. \quad (1)$$

We will now use the following lemma that shows that the surface area of any homogeneous polyhedral cone is independent of the number of faces k and in fact is at most 1 for all k .

Lemma 14 (Gaussian Surface Area of Homogeneous Cones [Naz03]). *Let C be a cone with apex at the origin (i.e., an intersection of arbitrarily many halfspaces all of whose boundaries contain the origin). Then C has Gaussian surface area $\Gamma(C)$ at most 1.*

Using Lemma 14 we obtain that $\sum_{i=1}^k \Gamma(F_i(\mathbf{V} \mathbf{x})) \leq 1$. Next, we observe that, when the positive numbers a_1, \dots, a_k satisfy $\sum_{i=1}^k a_i \leq 1$, it holds that $\sum_{i=1}^k a_i \sqrt{\log(1/a_i)} \leq \sqrt{\sum_{i=1}^k a_i \log(1/a_i)} \leq \sqrt{\log(k)}$ (using the fact that the uniform distribution maximizes the entropy). Using this fact and Equation (1), we obtain

$$\lim_{\theta \rightarrow 0} \frac{\Pr_{\mathbf{x} \sim \mathcal{N}_d} [C_1(\mathbf{x}) \neq C_2(\mathbf{x})]}{\theta} \leq c \sqrt{\log(k)}.$$

Thus, we have shown that, for sufficiently small θ , it holds that $\Pr_{\mathbf{x} \sim \mathcal{N}_d} [C_1(\mathbf{x}) \neq C_2(\mathbf{x})] \leq c \sqrt{\log(k)} \theta$, but, as we discussed in the start of the proof, the general bound follows directly from the bound for sufficiently small values of $\theta > 0$. \square

C Learning in Top- r Disagreement from Label Rankings

We prove the next result which corresponds to a proper learning algorithm for LSF in the presence of bounded noise with respect to the top- r disagreement.

Theorem 6 (Proper Top- r Learning Algorithm). *Fix $\eta \in [0, 1/2)$, $r \in [k]$ and $\epsilon, \delta \in (0, 1)$. Let \mathcal{D} be an η -noisy linear label ranking distribution satisfying the assumptions of Definition 2. There exists an algorithm that draws $N = \tilde{O}\left(\frac{d rk}{\epsilon(1-2\eta)^6} \log(1/\delta)\right)$ samples from \mathcal{D} , runs in $\text{poly}(N)$ time and, with probability at least $1 - \delta$, outputs a Linear Sorting function $h : \mathbb{R}^d \rightarrow \mathbb{S}_k$ that is ϵ -close in top- r disagreement to the target.*

The main result of this section is the next lemma, which directly implies the above theorem (using the same steps as the proof of Theorem 5).

Lemma 15 (Top- r Misclassification). *Let $r \in [k]$. Consider two matrices $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{k \times d}$ and let \mathcal{N}_d be the standard Gaussian in d dimensions. We have that*

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\sigma_{1..r}(\mathbf{U}\mathbf{x}) \neq \sigma_{1..r}(\mathbf{V}\mathbf{x})] \leq c \cdot k \cdot r \cdot \sqrt{\log(kr)} \cdot \max_{i \neq j} \theta(\mathbf{U}_i - \mathbf{U}_j, \mathbf{V}_i - \mathbf{V}_j),$$

where $c > 0$ is some universal constant.

Proof. Let us set $\sigma_{1..r}(\mathbf{W}\mathbf{x})$ denote the ordering of the top- r alternatives in the ranking $\sigma(\mathbf{W}\mathbf{x})$. Moreover, recall that $\sigma_\ell(\mathbf{W}\mathbf{x})$ denotes the alternative in the ℓ -th position of the ranking $\sigma(\mathbf{W}\mathbf{x})$. For two matrices $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{k \times d}$, we have that

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\sigma_{1..r}(\mathbf{U}\mathbf{x}) \neq \sigma_{1..r}(\mathbf{V}\mathbf{x})] = \sum_{j=1}^k \Pr_{\mathbf{x} \sim \mathcal{N}_d} \left[\bigcup_{\ell=1}^r \{j = \sigma_\ell(\mathbf{U}\mathbf{x}), j \neq \sigma_\ell(\mathbf{V}\mathbf{x})\} \right].$$

The first step is to understand the geometry of the set $\bigcup_{\ell=1}^r \{\mathbf{x} : j = \sigma_\ell(\mathbf{U}\mathbf{x})\} = \{\mathbf{x} : j \in \sigma_{1..r}(\mathbf{U}\mathbf{x})\}$ for $j \in [k]$. We have that this set is equal to

$$\mathcal{T}_{\mathbf{U}}^{(j)} = \bigcup_{S \subseteq [k]: |S| \leq r-1} \bigcap_{i \in S} \{\mathbf{x} : (\mathbf{U}_i - \mathbf{U}_j) \cdot \mathbf{x} \geq 0\} \cap \bigcap_{i \notin S} \{\mathbf{x} : (\mathbf{U}_i - \mathbf{U}_j) \cdot \mathbf{x} \leq 0\}.$$

In words, $\mathcal{T}_{\mathbf{U}}^{(j)}$ iterates over any possible collection of alternatives that can win the element j (they lie in the set of top elements S) and the remaining elements lose when compared with j (they lie in the complement set $[k] \setminus S$). Overloading the notation, let us define the mapping $T(\mathbf{t}) = T(t_1, \dots, t_k) = \sum_{S \subseteq [k]: |S| \leq r-1} \prod_{i \in S} \mathbf{1}\{t_i \geq 0\} \prod_{i \notin S} \mathbf{1}\{t_i \leq 0\}$. Using this mapping, we can define the indicator of the set $\mathcal{T}_{\mathbf{U}}^{(j)}$ as $T((\mathbf{U}_1 - \mathbf{U}_j) \cdot \mathbf{x}, \dots, (\mathbf{U}_k - \mathbf{U}_j) \cdot \mathbf{x})$. The top- r disagreement $\Pr_{\mathbf{x} \sim \mathcal{N}_d} [j \in \sigma_{1..r}(\mathbf{U}\mathbf{x}), j \notin \sigma_{1..r}(\mathbf{V}\mathbf{x})]$ is equal to:

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [T((\mathbf{U}_1 - \mathbf{U}_j) \cdot \mathbf{x}, \dots, (\mathbf{U}_k - \mathbf{U}_j) \cdot \mathbf{x}) = 1, T((\mathbf{V}_1 - \mathbf{V}_j) \cdot \mathbf{x}, \dots, (\mathbf{V}_k - \mathbf{V}_j) \cdot \mathbf{x}) = 0].$$

So we have that

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\sigma_{1..r}(\mathbf{U}\mathbf{x}) \neq \sigma_{1..r}(\mathbf{V}\mathbf{x})] = \sum_{j=1}^k \Pr_{\mathbf{x} \sim \mathcal{N}_d} [T_j(\mathbf{U}\mathbf{x}) = 1, T_j(\mathbf{V}\mathbf{x}) = 0] \leq \sum_{j=1}^k \Pr_{\mathbf{x} \sim \mathcal{N}_d} [T_j(\mathbf{U}\mathbf{x}) \neq T_j(\mathbf{V}\mathbf{x})].$$

In order to show the desired bound, it suffices to prove the following two lemmas.

Lemma 16 (Disagreement Region). *Consider a positive integer $r \leq k$. Fix $j \in [k]$ and let $\theta = \max_{i \in [k]} \theta(\mathbf{U}_i - \mathbf{U}_j, \mathbf{V}_i - \mathbf{V}_j)$. Then it holds that*

$$\lim_{\theta \rightarrow 0} \frac{\Pr_{\mathbf{x} \sim \mathcal{N}_d} [T_j(\mathbf{U}\mathbf{x}) \neq T_j(\mathbf{V}\mathbf{x})]}{\theta} \leq c \cdot \sum_{i \in [k]} \Gamma(F_i^j) \sqrt{\log \left(\frac{1}{\Gamma(F_i^j)} + 1 \right)},$$

where $c > 0$ is some constant and F_i^j is the surface $\{\mathbf{x} : j \in \sigma_{1..r}(\mathbf{V}\mathbf{x})\} \cap \{\mathbf{x} : \mathbf{V}_i \cdot \mathbf{x} = \mathbf{V}_j \cdot \mathbf{x}\}$ for the matrix $\mathbf{V} \in \mathbb{R}^{k \times d}$.

and,

Lemma 17. *Let F_i^j, r, k as in the previous lemma. It holds that*

$$\sum_{i \in [k]} \sum_{j \in [k]} \Gamma(F_i^j) \leq 2kr.$$

Applying these two lemmas with $\theta = \max_{i \neq j} \theta(\mathbf{U}_i - \mathbf{U}_j, \mathbf{V}_i - \mathbf{V}_j)$, we get that

$$Z := \lim_{\theta \rightarrow 0} \frac{\sum_{j \in [k]} \Pr_{\mathbf{x} \sim \mathcal{N}_d} [T_j(\mathbf{U}\mathbf{x}) \neq T_j(\mathbf{V}\mathbf{x})]}{\theta} \leq c \cdot \sum_{j \in [k]} \sum_{i \in [k]} \Gamma(F_i^j) \sqrt{\log \left(\frac{1}{\Gamma(F_i^j)} + 1 \right)}.$$

Let us set $\Gamma'(F_i^j) = \Gamma(F_i^j)/(2kr)$. Then we have that

$$Z \leq 2ckr \cdot \sum_{j \in [k]} \sum_{i \in [k]} \Gamma'(F_i^j) \sqrt{\log \left(\frac{1}{2kr \cdot \Gamma'(F_i^j)} + 1 \right)}.$$

It suffices to bound the quantity

$$\sum_{j \in [k]} \sum_{i \in [k]} \Gamma'(F_i^j) \sqrt{\log \left(\frac{1}{\Gamma'(F_i^j)} + 1 \right)} = O \left(kr \sqrt{\log(kr)} \right),$$

where we used a similar ‘‘entropy-like’’ inequality as we did in the top-1 case. This yields (by recalling that it is sufficient to consider only the case of arbitrarily small angles, as in the top-1 case) that

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [\sigma_{1..r}(\mathbf{U}\mathbf{x}) \neq \sigma_{1..r}(\mathbf{V}\mathbf{x})] \leq c r k \sqrt{\log(kr)} \cdot \max_{i \neq j} \theta(\mathbf{U}_i - \mathbf{U}_j, \mathbf{V}_i - \mathbf{V}_j),$$

for some universal constant c . \square

C.1 The proof of Lemma 16

We proceed with the proof of the key lemma concerning the disagreement region. We first show the following claim where we only change a single vector. Recall that

$$T(\mathbf{V}\mathbf{x}) = \sum_{S: |S| \leq r-1} \prod_{i \in S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \geq 0\} \prod_{i \notin S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \leq 0\}.$$

We will be interested in the surface $F_1 := F_1(\mathbf{V}\mathbf{x}) = T(\mathbf{V}\mathbf{x}) \mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} = 0\}$.

Claim 11. *Let $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^d$ and $\mathbf{r} \in \mathbb{R}^d$ with $\theta(\mathbf{v}_1, \mathbf{r}) \leq \theta$ for some sufficiently small $\theta \in (0, \pi/2)$. It holds that*

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [T(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \neq T(\mathbf{r} \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})] \leq c \cdot \theta \cdot \Gamma(F_1) \sqrt{\log \left(\frac{1}{\Gamma(F_1)} + 1 \right)},$$

where F_1 is the surface $T(\mathbf{V}\mathbf{x}) \cap \{\mathbf{x} : \mathbf{v}_1 \cdot \mathbf{x} = 0\}$ and c is some universal constant.

Proof. We first decompose the sum of $T(\mathbf{V}\mathbf{x})$ depending on whether $1 \in S$ or not. Hence, we have that $T(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) = T^+(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) + T^-(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})$ where

$$\begin{aligned} T^+(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) &= \sum_{S \subseteq [k]: |S| \leq r-1, 1 \in S} \prod_{i \in S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \geq 0\} \prod_{i \notin S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \leq 0\} \\ &= \sum_{S \subseteq [k]: |S| \leq r-1, 1 \in S} \mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \geq 0\} \cdot \prod_{i \in S \setminus \{1\}} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \geq 0\} \prod_{i \notin S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \leq 0\} \\ &= \mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \geq 0\} \cdot \sum_{S \subseteq [k]: |S| \leq r-1, 1 \in S} \prod_{i \in S \setminus \{1\}} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \geq 0\} \prod_{i \notin S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \leq 0\} \\ &=: \mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \geq 0\} \cdot G^+(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}), \end{aligned}$$

and similarly

$$\begin{aligned} T^-(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) &= \mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \leq 0\} \cdot \sum_{S \subseteq [k]: |S| \leq r-1, 1 \notin S} \prod_{i \in S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \geq 0\} \prod_{i \notin S \setminus \{1\}} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \leq 0\} \\ &=: \mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \leq 0\} \cdot G^-(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}). \end{aligned}$$

Notice that the indicator G^s does not depend on the alternative 1 for $s \in \{-, +\}$. Since $T : \mathbb{R}^k \rightarrow \{0, 1\}$, we have that

$$\begin{aligned} &\Pr_{\mathbf{x} \sim \mathcal{N}_d} [T(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \neq T(\mathbf{r} \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})] \\ &= \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [|T(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) - T(\mathbf{r} \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})|] \\ &\leq \sum_{s \in \{-, +\}} \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [|T^s(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) - T^s(\mathbf{r} \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})|] \\ &= \sum_{s \in \{-, +\}} \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [G^s(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \cdot |\mathbf{1}\{s \cdot \mathbf{v}_1 \cdot \mathbf{x} \geq 0\} - \mathbf{1}\{s \cdot \mathbf{r} \cdot \mathbf{x} \geq 0\}|]. \end{aligned}$$

Let us focus on the case $s = +$. The difference between the two indicators in the last line of the above equation corresponds to the event that the halfspaces $\mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \geq 0\}$ and $\mathbf{1}\{\mathbf{r} \cdot \mathbf{x} \geq 0\}$ disagree. Hence, we have that $|\mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} \geq 0\} - \mathbf{1}\{\mathbf{r} \cdot \mathbf{x} \geq 0\}| = \mathbf{1}\{(\mathbf{v}_1 \cdot \mathbf{x})(\mathbf{r} \cdot \mathbf{x}) < 0\}$. Note that the above indicator depends on both \mathbf{v}_1 and \mathbf{r} . We would like to work only with one of these two vectors. To this end, let us introduce \mathbf{q} , the normalized projection of \mathbf{r} onto the orthogonal complement of \mathbf{v}_1 , i.e., $\mathbf{q} = \text{proj}_{\mathbf{v}_1^\perp} \mathbf{r} / \|\text{proj}_{\mathbf{v}_1^\perp} \mathbf{r}\|_2$. We have that \mathbf{v}_1 and \mathbf{q} is an orthonormal basis of the subspace spanned by the vectors \mathbf{v}_1 and \mathbf{r} . Notice that $\mathbf{r} = \cos \theta(\mathbf{v}_1, \mathbf{r})\mathbf{v}_1 + \sin \theta(\mathbf{v}_1, \mathbf{r})\mathbf{q}$, by the construction of \mathbf{q} . Our goal is to understand the structure of the region $(\mathbf{v}_1 \cdot \mathbf{x})(\mathbf{r} \cdot \mathbf{x}) < 0$. This set is equal to

$$\{0 < \mathbf{v}_1 \cdot \mathbf{x} < -(\mathbf{q} \cdot \mathbf{x}) \tan \theta(\mathbf{v}_1, \mathbf{r})\} \cup \{-(\mathbf{q} \cdot \mathbf{x}) \tan \theta(\mathbf{v}_1, \mathbf{r}) < \mathbf{v}_1 \cdot \mathbf{x} < 0\}.$$

To see this, we have that $(\mathbf{v}_1 \cdot \mathbf{x})(\mathbf{r} \cdot \mathbf{x}) = (\mathbf{v}_1 \cdot \mathbf{x})(\cos \theta(\mathbf{v}_1, \mathbf{r})\mathbf{v}_1 \cdot \mathbf{x} + \sin \theta(\mathbf{v}_1, \mathbf{r})\mathbf{q} \cdot \mathbf{x})$. This quantity must be negative. The left-hand set considers the case where $\mathbf{v}_1 \cdot \mathbf{x} > 0$ and so $\tan \theta(\mathbf{v}_1, \mathbf{r})(\mathbf{q} \cdot \mathbf{x}) < -\mathbf{v}_1 \cdot \mathbf{x}$. We obtain the right-hand set in a similar way. Thus, we have that the disagreement region $(\mathbf{v}_1 \cdot \mathbf{x})(\mathbf{r} \cdot \mathbf{x}) < 0$ is a subset of the region $\{|\mathbf{v}_1 \cdot \mathbf{x}| \leq |\mathbf{q} \cdot \mathbf{x}| \tan \theta(\mathbf{v}_1, \mathbf{r})\}$. Since $\tan \theta(\mathbf{v}_1, \mathbf{r}) \leq \theta$ and we have that θ is sufficiently small we can also replace the above region by the larger region: $\{|\mathbf{v}_1 \cdot \mathbf{x}| \leq 2\theta|\mathbf{q} \cdot \mathbf{x}|\}$. Therefore, we have

$$\begin{aligned} & \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [G^+(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \mathbf{1}\{(\mathbf{v}_1 \cdot \mathbf{x})(\mathbf{r} \cdot \mathbf{x}) < 0\}] \\ & \leq \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [G^+(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \mathbf{1}\{|\mathbf{v}_1 \cdot \mathbf{x}| \leq 2\theta|\mathbf{q} \cdot \mathbf{x}|\}] . \end{aligned}$$

From this point, the proof goes as in the top-1 case. In total, we will get that

$$\begin{aligned} & \mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [T(\mathbf{v}_1 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) \neq T(\mathbf{r} \cdot \mathbf{x}, \mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})] \\ & = \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [(G^+(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x}) + G^-(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})) |\mathbf{q} \cdot \mathbf{x}| \delta(|\mathbf{v}_1 \cdot \mathbf{x}|)] \\ & \leq 2 \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| d\mu(\mathbf{x}) \\ & \leq 2 \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| \mathbf{1}\{|\mathbf{q} \cdot \mathbf{x}| \leq \xi\} d\mu(\mathbf{x}) + 2 \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| \mathbf{1}\{|\mathbf{q} \cdot \mathbf{x}| \geq \xi\} d\mu(\mathbf{x}) \\ & \leq 2\xi \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) d\mu(\mathbf{x}) + 2 \int_{\mathbf{x} \in F_1} \phi_d(\mathbf{x}) |\mathbf{q} \cdot \mathbf{x}| \mathbf{1}\{|\mathbf{q} \cdot \mathbf{x}| \geq \xi\} d\mu(\mathbf{x}), \end{aligned}$$

where $d\mu(\mathbf{x})$ is the standard surface measure in \mathbb{R}^d . Let us explain the first inequality above. Note that the space induced by $G^-(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})$ contains the space induced by $G^+(\mathbf{v}_2 \cdot \mathbf{x}, \dots, \mathbf{v}_k \cdot \mathbf{x})$. Hence, in the integration, we can integrate over the surface $F_1 = T(\mathbf{V}\mathbf{x}) \cap \mathbf{1}\{\mathbf{x} : \mathbf{v}_1 \cdot \mathbf{x} = 0\}$ twice. Essentially, this surface corresponds to $\mathbf{1}\{\mathbf{v}_1 \cdot \mathbf{x} = 0\} \cdot \sum_{S \subseteq [k] \setminus \{1\}; |S| \leq r-1} \prod_{i \in S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \geq 0\} \prod_{i \notin S} \mathbf{1}\{\mathbf{v}_i \cdot \mathbf{x} \leq 0\}$. Applying the steps of the top-1 case, we can obtain the desired bound in terms of the Gaussian surface area of F_1 . \square

Next, for fixed $j \in [k]$, we can apply the above claim sequentially (as we did in the end of the top-1 case) to get

$$\lim_{\theta \rightarrow 0} \frac{\mathbf{Pr}_{\mathbf{x} \sim \mathcal{N}_d} [T_j(\mathbf{U}\mathbf{x}) \neq T_j(\mathbf{V}\mathbf{x})]}{\theta} \leq c \cdot \sum_{i \in [k]} \Gamma(F_i^j) \sqrt{\log \left(\frac{1}{\Gamma(F_i^j)} + 1 \right)},$$

for some small constant $c > 0$.

C.2 The proof of Lemma 17

Using the above result, we get that it suffices to control the value $\Gamma(F_i^j)$, where F_i^j is the surface of $T_j(\mathbf{V}\mathbf{x}) \cap \{\mathbf{x} : \mathbf{V}_i \cdot \mathbf{x} = \mathbf{V}_j \cdot \mathbf{x}\}$ for the matrix \mathbf{V} and $i, j \in [k]$. We next have to control the Gaussian surface area of the induced shape, i.e., the quantity

$$\Gamma(\{\mathbf{x} : j \in \sigma_{1..r}(\mathbf{V}\mathbf{x})\} \cap \{\mathbf{x} : \mathbf{V}_i \cdot \mathbf{x} = \mathbf{V}_j \cdot \mathbf{x}\}).$$

To this end, we give the next lemma.

Lemma 18. Let $r \leq k$ with $r, k \in \mathbb{N}$. For any matrix $\mathbf{V} \in \mathbb{R}^{k \times d}$ and $i, j \in [k]$, there exists a matrix $\mathbf{Q} = \mathbf{Q}^{(i)} \in \mathbb{R}^{k \times d}$ which depends only on i such that

$$\Gamma(F_i^j) := \Gamma(\{\mathbf{x} : j \in \sigma_{1..r}(\mathbf{V}\mathbf{x})\} \cap \{\mathbf{x} : \mathbf{V}_i \cdot \mathbf{x} = \mathbf{V}_j \cdot \mathbf{x}\}) \leq 2 \cdot \Pr_{\mathbf{x} \sim \mathcal{N}_d} [j \in \sigma_{1..r}(\mathbf{Q}\mathbf{x})].$$

Before proving this result, let us see how to apply it in order to get Lemma 17. We will have that

$$\begin{aligned} \sum_{i \in [k]} \sum_{j \in [k]} \Gamma(F_i^j) &= \sum_{i \in [k]} \sum_{j \in [k]} \Gamma(\{\mathbf{x} : j \in \sigma_{1..r}(\mathbf{V}\mathbf{x})\} \cap \{\mathbf{x} : \mathbf{V}_i \cdot \mathbf{x} = \mathbf{V}_j \cdot \mathbf{x}\}) \\ &\leq 2 \sum_{i \in [k]} \sum_{j \in [k]} \Pr_{\mathbf{x} \sim \mathcal{N}_d} [j \in \sigma_{1..r}(\mathbf{Q}^{(i)}\mathbf{x})] \\ &= 2 \sum_{i \in [k]} \mathbf{E}_{\mathbf{x} \sim \mathcal{N}_d} [|\sigma_{1..r}(\mathbf{Q}^{(i)}\mathbf{x})|] \\ &= 2 \sum_{i \in [k]} r \\ &= 2kr. \end{aligned}$$

Proof of Lemma 18. For this proof, we fix $i, j \in [k]$. The first step is to design the matrix \mathbf{Q} . As a first observation, we can subtract the vector \mathbf{V}_i from each weight vector and do not affect the resulting orderings. Second, we can assume that the weight vectors that correspond to indices which j beats are unit. Let us be more specific Assume that initially we have that

$$(\mathbf{V}_j - \mathbf{V}_\ell) \cdot \mathbf{x} \geq 0.$$

The first observation gives that

$$(\mathbf{V}_j - \mathbf{V}_i) \cdot \mathbf{x} \geq (\mathbf{V}_\ell - \mathbf{V}_i) \cdot \mathbf{x}.$$

Let us set $\tilde{\mathbf{Q}}$ the intermediate matrix with rows $\mathbf{V}_j - \mathbf{V}_i$. The second observation states that the inequalities where j beats some index ℓ are not affected by normalization. Note that $\tilde{\mathbf{Q}}_j \cdot \mathbf{x} = 0$ and hence $\tilde{\mathbf{Q}}_\ell \cdot \mathbf{x} \leq 0$. Hence, dividing with non-negative numbers will not affect the order of these two values, i.e.,

$$\frac{\tilde{\mathbf{Q}}_j \cdot \mathbf{x}}{\|\tilde{\mathbf{Q}}_j\|_2} \geq \frac{\tilde{\mathbf{Q}}_\ell \cdot \mathbf{x}}{\|\tilde{\mathbf{Q}}_\ell\|_2}.$$

Note that the above ordering is \mathbf{x} -dependent, since the indices that j beats depend on \mathbf{x} . However, we can normalize any row of $\tilde{\mathbf{Q}}$ without affecting the fact that the element j is top- r (since the sign of the inner products is not affected by normalization). This transformation yields a matrix $\mathbf{Q} = \mathbf{Q}^{(i)}$ and depends only on i (crucially, it is independent of j). For simplicity, we will omit the index i in what follows. For this matrix, we have that

$$\{\mathbf{x} : j \in \sigma_{1..r}(\mathbf{Q}\mathbf{x}), \mathbf{Q}_j \cdot \mathbf{x} = 0\} = \{\mathbf{x} : j \in \sigma_{1..r}(\mathbf{V}\mathbf{x}), \mathbf{V}_i \cdot \mathbf{x} = \mathbf{V}_j \cdot \mathbf{x}\}.$$

We will now prove that

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [j \in \sigma_{1..r}(\mathbf{Q}\mathbf{x})] \geq \frac{\Gamma(F_i^j)}{2}.$$

Let us fix some \mathbf{x} and set $\mathbf{x}^\parallel = \text{proj}_{\mathbf{Q}_j} \mathbf{x}$ and $\mathbf{x}^\perp = \text{proj}_{\mathbf{Q}_j^\perp} \mathbf{x}$. We assume that \mathbf{x} lies in the set $\{\mathbf{x} : j \in \sigma_{1..r}(\mathbf{Q}\mathbf{x})\}$. This implies that there exist an index set I of size at least $k - r$ so that if $\ell \in I$ then

$$\mathbf{Q}_j \cdot \mathbf{x}^\parallel + \mathbf{Q}_j \cdot \mathbf{x}^\perp \geq \mathbf{Q}_\ell \cdot \mathbf{x}^\parallel + \mathbf{Q}_\ell \cdot \mathbf{x}^\perp.$$

Let us condition on the event

$$\mathbf{Q}_j \cdot \mathbf{x}^\perp \geq \mathbf{Q}_\ell \cdot \mathbf{x}^\perp.$$

We hence get that

$$\mathbf{Q}_j \cdot \mathbf{x}^\parallel = (\mathbf{Q}_j \cdot \mathbf{Q}_j) \cdot (\mathbf{Q}_j \cdot \mathbf{x}) \geq \mathbf{Q}_\ell \cdot \mathbf{x}^\parallel = (\mathbf{Q}_\ell \cdot \mathbf{Q}_j) \cdot (\mathbf{Q}_j \cdot \mathbf{x})$$

Using that \mathbf{Q}_j is unit, that the inner product between \mathbf{Q}_ℓ and \mathbf{Q}_j is at most one and that $\mathbf{Q}_j \cdot \mathbf{x}$ is a univariate Gaussian, we get that

$$\Pr_{z \sim \mathcal{N}(0,1)} [z \cdot (1 - \mathbf{Q}_\ell \cdot \mathbf{Q}_j) \geq 0] = 1/2.$$

The above discussion implies that

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [j \in \sigma_{1..r}(\mathbf{Q}\mathbf{x})] = \Pr_{\mathbf{x} \sim \mathcal{N}_d} [(\forall \ell \in I) \mathbf{Q}_j \cdot \mathbf{x}^\parallel + \mathbf{Q}_j \cdot \mathbf{x}^\perp \geq \mathbf{Q}_\ell \cdot \mathbf{x}^\parallel + \mathbf{Q}_\ell \cdot \mathbf{x}^\perp]$$

and so $\Pr_{\mathbf{x} \sim \mathcal{N}_d} [j \in \sigma_{1..r}(\mathbf{Q}\mathbf{x})]$ equals to

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [(\forall \ell \in I) \mathbf{Q}_j \cdot \mathbf{x}^\parallel \geq \mathbf{Q}_j \cdot \mathbf{x}^\parallel \mid (\forall \ell \in I) \mathbf{Q}_j \cdot \mathbf{x}^\perp \geq \mathbf{Q}_\ell \cdot \mathbf{x}^\perp] \cdot \Pr_{\mathbf{x} \sim \mathcal{N}_d} [(\forall \ell \in I) \mathbf{Q}_j \cdot \mathbf{x}^\perp \geq \mathbf{Q}_\ell \cdot \mathbf{x}^\perp].$$

However, in the above product, we have that the first term is 1/2 and the second term is the probability that $j \in \sigma_{1..r}(\mathbf{Q}\mathbf{x}^\perp)$, i.e.,

$$\Pr_{\mathbf{x} \sim \mathcal{N}_d} [j \in \sigma_{1..r}(\mathbf{Q}\mathbf{x})] \geq \frac{\Pr[j \in \sigma_{1..r}(\mathbf{Q}\mathbf{x}^\perp)]}{2} = \Gamma(F_i^j)/2,$$

since the space in the RHS is low-dimensional and corresponds to the desired surface. \square

D Distribution-Free Lower Bounds for Top-1 Disagreement Error

We begin with some definitions concerning the PAC Label Ranking setting. Let \mathcal{X} be an instance space and $\mathcal{Y} = \mathbb{S}_k$ be the space of labels, which are rankings over k elements. A sorting function or hypothesis is a mapping $h : \mathcal{X} \rightarrow \mathbb{S}_k$. We denote by $h_1(x)$ the top-1 element of the ranking $h(x)$. A hypothesis class is a set of classifiers $\mathcal{H} \subset \mathbb{S}_k^{\mathcal{X}}$.

Top-1 Disagreement Error. The top-1 disagreement error with respect to a joint distribution \mathcal{D} over $\mathcal{X} \times \mathbb{S}_k$ equals to the probability $\Pr_{(x,\sigma) \sim \mathcal{D}} [h_1(x) \neq \sigma^{-1}(1)]$. We mainly consider learning in the **realizable** case, which means that there is $h^* \in \mathcal{H}$ which has (almost surely) zero error. Therefore, we can focus on the marginal distribution \mathcal{D}_x over \mathcal{X} and denote the top-1 disagreement error of a sorting function h with respect to the true hypothesis h^* by $\text{Err}_{\mathcal{D}_x, h^*}(h) := \Pr_{x \sim \mathcal{D}_x} [h_1(x) \neq h_1^*(x)]$.

A learning algorithm is a function \mathcal{A} that receives a training set of m instances, $S \in \mathcal{X}^m$, together with their labels according to h^* . We denote the restriction of h^* to the instances in S by $h^*|_S$. The output of the algorithm \mathcal{A} , denoted $\mathcal{A}(S, h^*|_S)$ is a sorting function. A learning algorithm is proper if it always outputs a hypothesis from \mathcal{H} .

The top-1 PAC Label Ranking sample complexity of a learning algorithm \mathcal{A} is the function $m_{\mathcal{A}, \mathcal{H}}^{(1)}$ defined as follows: for every $\epsilon, \delta > 0$, $m_{\mathcal{A}, \mathcal{H}}^{(1)}(\epsilon, \delta)$ is the minimal integer such that for every $m \geq m_{\mathcal{A}, \mathcal{H}}^{(1)}(\epsilon, \delta)$, every distribution \mathcal{D}_x on \mathcal{X} , and every target hypothesis $h^* \in \mathcal{H}$, $\Pr_{S \sim \mathcal{D}_x^m} [\text{Err}_{\mathcal{D}_x, h^*}(\mathcal{A}(S, h^*|_S)) > \epsilon] \leq \delta$. In this case, we say that the learning algorithm (ϵ, δ) -learns the class of sorting functions \mathcal{H} with respect to the top-1 disagreement error. If no integer satisfies the inequality above, define $m_{\mathcal{A}}^{(1)}(\epsilon, \delta) = \infty$. \mathcal{H} is learnable with \mathcal{A} if for all ϵ and δ the sample complexity is finite. The **top-1 PAC Label Ranking sample complexity** of a class \mathcal{H} is $m_{\text{PAC}, \mathcal{H}}^{(1)}(\epsilon, \delta) = \inf_{\mathcal{A}} m_{\mathcal{A}, \mathcal{H}}^{(1)}(\epsilon, \delta)$, where the infimum is taken over all learning algorithms. Clearly, the above top-1 definition can be extended to the top- r setting.

In this section, we show the next result. We denote by $\mathcal{L}_{d,k}$ the class of Linear Sorting functions in d dimensions with k labels.

Theorem 7. *In the realizable PAC Label Ranking setting, any algorithm that (ϵ, δ) -learns the class $\mathcal{L}_{d,k}$ with respect to the top-1 disagreement error requires at least $\Omega((dk + \log(1/\delta))/\epsilon)$ samples.*

D.1 Top-1 Ranking Natarajan Dimension

In order to establish the above result, we introduce a variant of the standard Natarajan dimension [Nat89, BDCBL92, DSBDS11, DSS14]. For a ranking π , we will also let $L_1(\pi)$ its top-1 element and $L_{3..k}(\pi)$ the ranking after deleting its top-2 part.

Definition 3 (Top-1 Ranking Natarajan Dimension). Let $\mathcal{H} \subseteq \mathbb{S}_k^{\mathcal{X}}$ be a hypothesis class of sorting functions and let $S \subseteq \mathcal{X}$. We say that \mathcal{H} N -shatters S if there exist two mappings $f_1, f_2 : S \rightarrow \mathbb{S}_k$ such that for every $y \in S$, $L_1(f_1(y)) \neq L_1(f_2(y))$ and $L_{3..k}(f_1(y)) = L_{3..k}(f_2(y))$ and for every $T \subseteq S$, there exists a sorting function $g \in \mathcal{H}$ such that

$$(i) \forall x \in T, g(x) = f_1(x), \text{ and } (ii) \forall x \in S \setminus T, g(x) = f_2(x).$$

The **top-1 Ranking Natarajan dimension** of \mathcal{H} , denoted $d_N^{(1)}(\mathcal{H})$ is the maximal cardinality of a set that is N -shattered by \mathcal{H} .

First, we connect PAC Label Ranking learnability to the top-1 disagreement error with the notion of top-1 Ranking Natarajan dimension.

Theorem 8 (Top-1-Natarajan Lower Bounds Sample Complexity). In the realizable PAC Label Ranking setting, we have for every hypothesis class $\mathcal{H} \subseteq \mathbb{S}_k^{\mathcal{X}}$

$$m_{\text{PAC}, \mathcal{H}}^{(1)}(\epsilon, \delta) = \Omega\left(\frac{d_N^{(1)}(\mathcal{H}) + \ln(1/\delta)}{\epsilon}\right).$$

Proof. Let $\mathcal{H} \subseteq \mathbb{S}_k^{\mathcal{X}}$ be a class of sorting functions of top-1-Natarajan dimension $d_N^{(1)} = d_N$. Consider the binary hypothesis class $\mathcal{H}_{\text{bin}} = \{0, 1\}^{[d_N]}$ which contains all the classifiers from $[d_N] = \{1, \dots, d_N\}$ to $\{0, 1\}$. It suffices to show the following.

Claim 12. It holds that $m_{\text{PAC}, \mathcal{H}}^{(1)}(\epsilon, \delta) \geq m_{\text{PAC}, \mathcal{H}_{\text{bin}}}(\epsilon, \delta)$.

This is sufficient since we have that $m_{\text{PAC}, \mathcal{H}_{\text{bin}}}(\epsilon, \delta) = \Omega\left(\frac{\text{VC}(\mathcal{H}_{\text{bin}}) + \ln(1/\delta)}{\epsilon}\right)$ and $\text{VC}(\mathcal{H}_{\text{bin}}) = d_N$. Let us now prove the claim.

We assume that the instance space is the set \mathcal{X} . Assume that A is a learning algorithm for the hypothesis class $\mathcal{H} \subseteq \mathbb{S}_k^{\mathcal{X}}$ and A_{bin} is a learning algorithm for the associated binary class \mathcal{H}_{bin} . It suffices to show that A requires at least as many samples as A_{bin} . In fact, we will show that whenever A_{bin} errs, so does A . Let $S = \{s_1, \dots, s_{d_N}\}$, f_0, f_1 be the set and the two functions that witness that the top-1-Natarajan dimension of \mathcal{H} is d_N . Given a training set $(x_i, y_i)_{i \in [m]} \in ([d_N] \times \{0, 1\})^m$, we set $g : \mathcal{X} \rightarrow \mathbb{S}_k$ be equal to the output of the algorithm A with input $(s_{x_i}, f_{y_i}(x_i))_{i \in [m]} \in (S \times \mathbb{S}_k)^m$. We also set f be the output of the algorithm A_{bin} with input $(x_i, y_i)_{i \in [m]}$ by setting $f(i) = 1$ if and only if $L_1(g(s_i)) = L_1(f_1(s_i))$. We will show that whenever A_{bin} errs, so does A . Fix $(x_i, y_i) \in S \times \{0, 1\}$. Assume that $A_{\text{bin}}(x_i) \neq y_i$ and say $y_i = 0$. Then $f(i) = 1$ and so $L_1(g(s_i)) = L_1(f_1(s_i)) \neq L_1(f_0(s_i))$. This implies that A errs. The case $y_i = 1$ is similar. \square

D.2 Lower Bound for top-1 disagreement error for LSFs

Theorem 9 (Top-1 Natarajan Dimension of LSFs). Consider the hypothesis class $\mathcal{L}_{d,k} = \{\sigma_{\mathbf{W}} : \mathbb{R}^d \rightarrow \mathbb{S}_k : \sigma_{\mathbf{W}}(\mathbf{x}) = \text{argsort}(\mathbf{W}\mathbf{x}), \mathbf{W} \in \mathbb{R}^{k \times d}\}$. Then, $d_N^{(1)}(\mathcal{L}_{d,k}) = \Omega(dk)$.

Proof. Fix $k \in \mathbb{N}$. Let us consider the case $d = 2$ that will correspond as the building block for the general case $d > 2$. Let us first choose the set of points: Set P be the collection of pairs $P = \{(2i-1, 2i)\}_{i \in [b]}$ for any $i \in [b]$ with $b = \lfloor k/2 \rfloor$ and $S = \{\mathbf{x}_m\}_{m \in P}$ where these points correspond to $|P|$ equidistributed points on the unit sphere in \mathbb{R}^2 . This set of points has size $|P| = \Theta(k)$ and we are going to N -shatter it using $\mathcal{L}_{2,k}$.

Consider the matrix $\mathbf{W} \in \mathbb{R}^{k \times 2}$ so that $\{\mathbf{W}_i\}_{i \in [k]}$ correspond to the rows of \mathbf{W} . The structure of the problem relies on the hyperplanes with normal vectors $(\mathbf{W}_i - \mathbf{W}_j)_{i \neq j}$ and our choice of \mathbf{W} will rely on these hyperplanes. For any $m = (2i-1, 2i)$, we set $\mathbf{W}_{2i-1}, \mathbf{W}_{2i}$ on the unit sphere so that $\mathbf{W}_{2i-1} \cdot \mathbf{W}_{2i} = 1 - \phi$ with $\phi \in (0, 1)$ sufficiently small (set $\arccos(1 - \phi) = 2\pi/(100k)$) and let C_m be the cone generated by these two vectors with axis I_m . We place \mathbf{W}_{2i-1} so that the distance between \mathbf{x}_m and the hyperplane I_m is sufficiently small (say that the angle between \mathbf{x}_m and I_m is $\arccos(1 - \phi)/100$). Note that the normal vector of I_m is $\mathbf{W}_{2i-1} - \mathbf{W}_{2i}$ and we place \mathbf{x}_m so that it has positive correlation with this vector. This uniquely identifies the location of \mathbf{W}_{2i} . Crucially, each vector \mathbf{x}_m has the following properties: (i) \mathbf{x}_m is very close to the boundary of the hyperplane

with normal vector $(\mathbf{W}_{2i-1} - \mathbf{W}_{2i})$, (ii) $\mathbf{W}_{2i-1} \cdot \mathbf{x}_m > \mathbf{W}_{2i} \cdot \mathbf{x}_m > \mathbf{W}_j \cdot \mathbf{x}_m$ for any $j \notin m$ and (iii) \mathbf{x}_m is far from any boundary induced by hyperplanes with normal vectors $\mathbf{W}_j - \mathbf{W}_{j'}$ for any $(j, j') \neq m$.

Since the points are well-separated on the unit sphere, for any $m = (2i - 1, 2i) \in P$, we have $\mathbf{W}_{2i-1} \cdot \mathbf{W}_{2i} = 1 - \phi \approx 1$ and for any other pair of indices $(i, j) \notin P$, there exists $c = c(k) \in (0, 1)$, $|\langle \mathbf{W}_i, \mathbf{W}_j \rangle| \leq c$.

For any $m = (2i - 1, 2i) \in P$, we set $\mathbf{W}'_{2i-1} - \mathbf{W}'_{2i} = \mathbf{R}_\theta(\mathbf{W}_{2i-1} - \mathbf{W}_{2i})$ for some θ to be chosen, where \mathbf{R}_θ is the 2×2 rotation matrix. We choose θ so that each point \mathbf{x}_m for $m = (2i - 1, 2i) \in P$ with $(\mathbf{W}_{2i-1} - \mathbf{W}_{2i}) \cdot \mathbf{x}_m > 0$ satisfies $(\mathbf{W}'_{2i-1} - \mathbf{W}'_{2i}) \cdot \mathbf{x}_m < 0$. The main idea is that since \mathbf{x}_m has the properties (i)-(iii) described above, the rankings induced by the vectors $\mathbf{W}\mathbf{x}_m$ and $\mathbf{W}'\mathbf{x}_m$ will be different in the first two positions but the same in the rest.

Given the training set $\{\mathbf{x}_m\}_{m \in P}$, we have to construct f_0, f_1 and verify that they satisfy the top-1 Ranking Natarajan conditions. For $m = (2i - 1, 2i)$, we have that $f_0(\mathbf{x}_m) = (2i - 1, 2i, \pi)$ and $f_1(\mathbf{x}_m) = (2i, 2i - 1, \pi)$ for some ranking π of size $k - 2$ that depends on m . Specifically, we will set $f_0(\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x})$ and $f_1(\mathbf{x}) = \sigma(\mathbf{W}'\mathbf{x})$, where σ gives the decreasing ordering of the elements of the input vector. By the choice of the set S and \mathbf{W}, \mathbf{W}' , it remains to show that the $k - 2$ last elements of the rankings $f_0(\mathbf{x}_m)$ (say π_0) and of $f_1(\mathbf{x}_m)$ (say π_1) are in the same order, i.e., $L_{3..k}(f_0(\mathbf{x}_m)) = L_{3..k}(f_1(\mathbf{x}_m))$. Assume that $u \succ v$ in π_0 . It suffices to show that $(\mathbf{W}'_u - \mathbf{W}'_v) \cdot \mathbf{x}_m \geq 0$, i.e., the order of u and v is preserved when transforming \mathbf{W} to \mathbf{W}' . We have that $(\mathbf{W}_u - \mathbf{W}_v) \cdot \mathbf{x}_m > c_1$ for some constant $c_1 > 0$ (c_1 is the minimum over $(u, v) \neq m = (2i - 1, 2i)$). Hence, we can pick θ small enough so that $(\mathbf{W}'_u - \mathbf{W}'_v) \cdot \mathbf{x}_m > c_2$ and this can be done for any pair u, v that does not correspond to m . This implies that $\pi_0 = \pi_1 = \pi$. In particular, we have that

$$(\mathbf{W}'_u - \mathbf{W}'_v) \cdot \mathbf{x}_m = \cos(\theta) \cdot (\mathbf{W}_u - \mathbf{W}_v) \cdot \mathbf{x}_m + \sin(\theta) \cdot (W_{uv}^{(1)}x_m^{(2)} - W_{uv}^{(2)}x_m^{(1)}) > c_2 > 0$$

for some θ sufficiently small, where $W_{uv}^{(t)}$ is the t -th entry of the vector $\mathbf{W}_u - \mathbf{W}_v$ for $t \in \{1, 2\}$ and $\mathbf{x}_m, \mathbf{W}_u, \mathbf{W}_v$ are unit vectors.

For any subset T of S , it remains to choose a linear classifier in $\mathcal{L}_{2,k}$ (which is allowed to depend on T). For any $T \subseteq S = \{\mathbf{x}_m\}_{m \in P}$, we consider the matrix $\overline{\mathbf{W}} \in \mathbb{R}^{k \times 2}$ so that for the i -th row $\overline{\mathbf{W}}_i = \mathbf{W}_i \mathbf{1}\{i \in m \in T\} + \mathbf{W}'_i \mathbf{1}\{i \in m \in S \setminus T\}$ for any $i \in [k]$. This is valid since the pairs $m \in P$ partition $[k]$. We have to show the following two properties: (i) $\sigma(\overline{\mathbf{W}}\mathbf{x}) = f_0(\mathbf{x})$ for $\mathbf{x} \in T$ and (ii) $\sigma(\overline{\mathbf{W}}\mathbf{x}) = f_1(\mathbf{x})$ for $\mathbf{x} \in S \setminus T$.

Assume that $m = (2i - 1, 2i)$ and $\mathbf{x}_m \in T$. We have that $f_0(\mathbf{x}_m) = (2i - 1, 2i, \pi)$ and $\overline{\mathbf{W}}_{2i-1} - \overline{\mathbf{W}}_{2i} = \mathbf{W}_{2i-1} - \mathbf{W}_{2i}$ and so $2i - 1 \succ 2i$ in the ranking $\sigma(\overline{\mathbf{W}}\mathbf{x}_m)$. It remains to show that the remaining $\binom{k}{2} - 1$ pairwise comparisons are the same in the two rankings. Let us consider a pair of points $u \neq v$ so that $u \succ v$ in $f_0(\mathbf{x}_m)$. It suffices to show that $u \succ v$ in $\sigma(\overline{\mathbf{W}}\mathbf{x}_m)$.

1. If u, v are so that $\overline{\mathbf{W}}_u - \overline{\mathbf{W}}_v = \mathbf{W}_u - \mathbf{W}_v$, the result holds.
2. If u, v are so that $\overline{\mathbf{W}}_u - \overline{\mathbf{W}}_v = \mathbf{W}_u - \mathbf{W}'_v$: In this case, u and v lie in a different pair of P and this implies that the correct direction is preserved if θ is appropriately chosen. For θ as above, it holds that $(\mathbf{W}_u - \mathbf{R}_\theta \mathbf{W}_v) \cdot \mathbf{x}_m$ has the same sign as $(\mathbf{W}_u - \mathbf{W}_v) \cdot \mathbf{x}_m$. In particular,

$$\mathbf{W}_u \cdot \mathbf{x}_m - \mathbf{R}_\theta \mathbf{W}_v \cdot \mathbf{x}_m = \mathbf{W}_u \cdot \mathbf{x}_m - (\cos(\theta)W_v^{(1)} - \sin(\theta)W_v^{(2)})x_m^{(1)} - (\sin(\theta)W_v^{(1)} + \cos(\theta)W_v^{(2)})x_m^{(2)},$$

and so

$$(\mathbf{W}_u - \mathbf{W}'_v) \cdot \mathbf{x}_m = \cos(\theta) \cdot (\mathbf{W}_u - \mathbf{W}_v) \cdot \mathbf{x}_m + \sin(\theta)(W_v^{(2)}x_m^{(1)} - W_v^{(1)}x_m^{(2)}) > 0.$$

3. If u, v are so that $\overline{\mathbf{W}}_u - \overline{\mathbf{W}}_v = \mathbf{W}'_u - \mathbf{W}'_v$, the analysis for the inner product with \mathbf{x}_m will be similar.

We now have to extend this proof for $d > 2$. We will “tensorize” the above construction as follows. Let $S = \{\mathbf{y}_{mj}\}_{m \in [b], j \in [d/2]}$ with $|S| = \lfloor k/2 \rfloor \cdot \lfloor d/2 \rfloor$. We first define the points of S : For $s \in [d]$,

set $y_{m,j}[s] = x_m[1]1\{s = 2j - 1\} + x_m[2]1\{s = 2j\}$ with $\mathbf{y}_{m,j} \in \mathbb{R}^d$, i.e., $\mathbf{y}_{m,j}$ has the values of \mathbf{x}_m at the consecutive entries indicated by $m = (2i - 1, 2i) \in P$ and zeros at the other positions.

We have to show that the set S is N -shattered. Given $T \subseteq S$, we are going to create the matrix $\overline{\mathbf{W}} \in \mathbb{R}^{k \times d}$. For illustration, think of each row of the matrix as having $d/2$ blocks of size two. If $\mathbf{y}_{m,j} \in T$ with $m = (2i - 1, 2i)$, set the two associated rows (indicated by m) of $\overline{\mathbf{W}}$ with $\mathbf{W}_{2i-1}, \mathbf{W}_{2i}$ at the j -th block and with $\mathbf{W}'_{2i-1}, \mathbf{W}'_{2i}$ otherwise. We will have that $\sigma(\overline{\mathbf{W}}\mathbf{y}) = f_0(\mathbf{y})$ if $\mathbf{y} \in T$ and $\sigma(\overline{\mathbf{W}}\mathbf{y}) = f_1(\mathbf{y})$ otherwise and the analysis is the same as the $d = 2$ case. \square

E Examples of Noisy Ranking Distributions

Definition 4 (Mallows model [Mal57]). *Consider k alternatives and let $\pi \in \mathbb{S}_k, \phi \in [0, 1]$. The Mallows distribution $\mathcal{M}_{\text{Mal}}(\pi, \phi)$ with central ranking π and spread parameter ϕ is a probability measure over \mathbb{S}_k with density $\Pr_{\sigma \sim \mathcal{M}_{\text{Mal}}(\pi, \phi)}[\sigma]$ that is proportional to $\phi^{d(\sigma, \pi)}$, where d is a ranking distance.*

We focus on Mallows models associated with the Kendall's Tau distance $d = d_{KT}$ (the standard distance, not the normalized one), which measures the number of discordant pairs.

Fact 2. *When $\phi < 1$, the Mallows model $\mathcal{M}_{\text{Mal}}(\pi, \phi)$ is a ranking distribution with bounded noise at most $\frac{1+\phi}{4} < 1/2$.*

Proof. The following property holds [Mal57]

$$\Pr_{\sigma \sim \mathcal{M}_{\text{Mal}}(\pi, \phi)}[\sigma(i) < \sigma(j) | \pi(i) < \pi(j)] = \frac{\pi(j) - \pi(i) + 1}{1 - \phi^{\pi(j) - \pi(i) + 1}} - \frac{\pi(j) - \pi(i)}{1 - \phi^{\pi(j) - \pi(i)}} \geq \frac{1}{2} + \frac{1 - \phi}{4}.$$

\square

The Bradley-Terry-Luce model [BT52, Luc12] is the most studied pairwise comparisons model. In his seminal paper, Mallows [Mal57] also studied the following natural ranking distribution:

Definition 5 (Bradley-Terry-Mallows [Mal57]). *Consider a score vector $\mathbf{w} \in \mathbb{R}_+^k$ with k distinct entries and let π be the ranking induced by the values of \mathbf{w} in decreasing order. The Bradley-Terry-Mallows distribution $\mathcal{M}_{\text{BTM}}(\mathbf{w})$ with central ranking π is a probability measure over \mathbb{S}_k with density $\Pr_{\sigma \sim \mathcal{M}_{\text{BTM}}(\mathbf{w})}[\sigma]$ that is proportional to $\prod_{i \succ_{\sigma} j} \frac{w_i}{w_i + w_j}$.*

Lemma 19. *There exists a real number $0 < \eta < 1/2$ so that the Bradley-Terry-Mallows distribution $\mathcal{M}_{\text{BTM}}(\mathbf{w})$ is a ranking distribution with bounded noise at most η .*

Proof. In the standard Bradley-Terry-Luce model, the pairwise comparison between the alternatives i, j is a Bernoulli random variable with $\Pr[i \succ j] = w_i / (w_i + w_j)$. The Bradley-Terry-Mallows distribution can be considered as the Bradley-Terry-Luce model conditioned on the event that all the pairwise comparisons are consistent to a ranking. Hence, we have that

$$\Pr_{\sigma \sim \mathcal{M}_{\text{BTM}}(\mathbf{w})}[\sigma] = \frac{1}{Z(k, \mathbf{w})} \prod_{i \succ_{\sigma} j} \frac{w_i}{w_i + w_j}.$$

Let us set $\mathcal{A}_{i \succ j} = \{\sigma \in \mathbb{S}_k : \sigma(i) < \sigma(j)\}$. We are interested in the following probability

$$\Pr_{\sigma \sim \mathcal{M}_{\text{BTM}}(\mathbf{w})}[i \succ_{\sigma} j | w_i > w_j] = \Pr_{\sigma \sim \mathcal{M}_{\text{BTM}}(\mathbf{w})}[\sigma(i) < \sigma(j) | w_i > w_j] = \frac{1}{Z(k, \mathbf{w})} \sum_{\sigma \in \mathcal{A}_{i \succ j}} \prod_{p \succ_{\sigma} q} \frac{w_p}{w_p + w_q}.$$

Note that in order to show the desired property, it suffices to show that

$$\sum_{\sigma \in \mathcal{A}_{i \succ j}} \prod_{p \succ_{\sigma} q} \frac{w_p}{w_p + w_q} > \sum_{\sigma \in \mathcal{A}_{i \prec j}} \prod_{p \succ_{\sigma} q} \frac{w_p}{w_p + w_q}.$$

First, observe that there exists a correspondence mapping $\sigma \in \mathcal{A}_{i \succ j}$ to $\mathcal{A}_{i \prec j}$, where one flips the elements i and j . Hence, it suffices to show that the mass of the ranking $(u_a)i(u_b)j(u_c)$ is larger than the one of the ranking $(u_a)j(u_b)i(u_c)$, where u_a, u_b, u_c are permutations of length between 0 and

$k - 2$ with elements in $[k] \setminus \{i, j\}$. For the two above rankings, the only terms of the product that are not identical are the following

$$\frac{w_i}{w_i + w_j} \prod_{x \in u_b} \frac{w_i}{w_i + w_x} \frac{w_x}{w_x + w_j} > \frac{w_j}{w_i + w_j} \prod_{x \in u_b} \frac{w_j}{w_j + w_x} \frac{w_x}{w_x + w_i},$$

since $w_i > w_j$ and so the result follows. □