

SOLVING RANDOM SYSTEMS OF QUADRATIC EQUATIONS WITH TANH WIRTINGER FLOW

ZHENWEI LUO, YE ZHANG

Rice University
Houston, Texas

1. INTRODUCTION

Reconstructing signal from intensity measurements only, which is also known as phase retrieval, is an important problem with applications in various fields, including the X-ray crystallography, astronomy, and diffraction imaging [6, 8, 9]. There is a recent resurgence of interest in solving phase retrieval problem in machine learning community. A typical setting of those works is that, the observed data y_i is of the form,

$$y_i \approx |\langle \mathbf{a}_i, \mathbf{x} \rangle|^2, \quad i = 1, \dots, m,$$

where \mathbf{a}_i is a random gaussian measurement vector and \mathbf{x} is unknown [4]. For simplicity, we consider only the real-valued case in this paper and focus on the wirtinger flow type algorithms. The Wirtinger Flow (WF) method was first proposed by Candes et al. to solve the phasing problem [3]. Latterly, Chen et al. improved WF by wisely discarding certain outlier gradients, and named this procedure as the Truncated Wirtinger Flow (TWF) [5]. There are numerous follow up works focus on improving the truncation rules [10, 18, 19]. In this paper, we propose a new way to improve the success rate of solving phasing problem. Instead of designing new truncation rules, our method relies on a nonlinear activation function—tanh. This new method can not only significantly improve the success rate of solving the random systems of quadratic equations with low measurements to unknowns ratio, but also greatly increase the convergence rate.

2. TANH WIRTINGER FLOW

In this paper, we denote the true real signal as $\mathbf{x} \in \mathbb{R}^n$, and the design matrix as $\mathbf{A} \equiv [\mathbf{a}_1, \dots, \mathbf{a}_m]^T \in \mathbb{R}^{m \times n}$, where $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. We also use $O(n)$ to represent a quantity which is of the order Cn , where C is a constant grater than 1. Meanwhile, $o(n)$ denotes a quantity which is of the order cn , where c is a constant smaller than 1. Given an estimated signal \mathbf{z} and let $\mathbf{h} = \mathbf{x} - \mathbf{z}$, we then have the relationship between the observed data and the estimated data,

$$y_i \approx |\mathbf{a}_i^T \mathbf{z} + \mathbf{a}_i^T \mathbf{h}|^2 = |\mathbf{a}_i^T \mathbf{z}|^2 + 2\mathbf{a}_i^T \mathbf{z} \mathbf{a}_i^T \mathbf{h} + |\mathbf{a}_i^T \mathbf{h}|^2.$$

Since $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\mathbf{a}_i^T \mathbf{h}$ is a random variable and distributed according to $\mathcal{N}(0, |\mathbf{h}|^2 = \sigma^2)$. We first use a likelihood heuristic to derive the new optimization method. Suppose the additive noise of the observation y_i is negligible comparing with the noise generated by model errors and abbreviate the $\mathbf{a}_i^T \mathbf{h}$ and $\mathbf{a}_i^T \mathbf{z}$ as ϵ_i and f_i , we have

$$(1) \quad y_i = f_i^2 + 2f_i\epsilon_i + \epsilon_i^2.$$

Hence, we can derive the probability density function of the observed data conditioned on a given solution as following. From Equation 1, the noise can take two values for given y_i and f_i ,

$$\epsilon_i = -f_i \pm \sqrt{y_i}.$$

Moreover, since ϵ is distributed according to $\mathcal{N}(0, \sigma^2)$, the probability density function of observed data y_i conditioned on f_i is,

$$p(y_i|f_i) = \sum p(\epsilon|f_i) \left| \frac{d\epsilon}{dy_i} \right| = \frac{1}{\sqrt{2\pi y_i} \sigma} \exp\left(-\frac{f_i^2 + y_i^2}{2\sigma^2}\right) \cosh\left(\frac{f_i \sqrt{y_i}}{\sigma^2}\right),$$

which gives a likelihood function to estimate \mathbf{x} . We can obtain an estimation \mathbf{z} for \mathbf{x} by maximizing the total likelihood $\prod_{i=1}^m l_i(\mathbf{z}; y_i)$. The origin of this likelihood function can be traced to the development of maximum likelihood refinement method in crystallography fields [15, 12, 2]. General overviews are given by Murshudov et al. and Lunin et al. [12, 11]. Therefore, ignoring the terms consists only of observations, the random systems of quadratic equations can be solved by minimizing the following target function,

$$(2) \quad \min_{\mathbf{x}} \frac{1}{2m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{z}|^2 - 2\sigma^2 \log(\cosh \frac{\mathbf{a}_i^T \mathbf{z} \sqrt{y_i}}{\sigma^2}).$$

The gradient of target function 2 for any $\mathbf{z} \in \mathbb{R}^n$ is of the form,

$$(3) \quad \frac{1}{2m} \sum_{i=1}^m \nabla l_i(\mathbf{z}) = \frac{1}{m} \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{z} - \sqrt{y_i} \tanh \frac{\mathbf{a}_i^T \mathbf{z} \sqrt{y_i}}{\sigma^2}) \mathbf{a}_i.$$

This gives rise to a simple gradient descent update rule

$$\mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu}{2m} \sum_{i=1}^m \nabla l_i(\mathbf{z}).$$

The vanilla gradient descent is not the only method to update parameters. In fact, we can incorporate our new gradient with any first order optimization algorithm. To gain an empirical understanding about the new gradient, we consider it in the noiseless setting, namely, $y_i = |\mathbf{a}_i^T \mathbf{x}|^2$. Thus we can rewrite the gradient in equation 3 as

$$\frac{1}{2m} \sum_{i=1}^m \nabla l_i(\mathbf{z}) = \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^T (\mathbf{z} - \mathbf{x} \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}),$$

where $\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^T \mathbf{x} \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}$ can be viewed as an estimator for the true signal \mathbf{x} . Hence, the true signal is a linear combination of measurement vectors \mathbf{a}_i where each vector is weighted by $\mathbf{a}_i^T \mathbf{x} \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}$. It is worth noting that the square norm of error, σ^2 , should be known before computing the gradient. However, in practice, we don't have oracular knowledge about this factor. Thus we should design effective methods to estimate the square norm of error. However, a geometric observation about phase retrieval problem saves us from the problem of error estimation. According to the Grothendieck's identity from the Lemma 3.6.6 in [17], the relationship between the inner products of a random gaussian vector $\mathbf{a}_i \in \mathbb{R}^n$ with any fixed vectors $\mathbf{x}, \mathbf{z} \in S^{n-1}$ can be understood by reducing the problem to \mathbb{R}^2 . Using the reduction approach given in section 4.3, we have $\mathbf{a}_i^T \mathbf{x} = r \|\mathbf{x}\| \cos \theta$, $\mathbf{a}_i^T \mathbf{z} = r \|\mathbf{z}\| \cos(\theta - \phi)$. If we can design a weighting function which assigns smaller weight to the observation and estimation pairs with opposite

sign, our estimator will be closer to the true signal \mathbf{x} . In fact, $\cos \theta \cos(\theta - \phi)$ is an ideal candidate for such weighting functions since its average value are smaller in the region where $\text{sign}(\cos \theta) \text{sign}(\cos(\theta - \phi))$ is negative. That suggests we should estimate the length of projected design vector. As the square norm of projected design vector is given by $r^2 = \frac{1}{1 - \cos^2 \theta} (\frac{(\mathbf{a}_i^T \mathbf{z})^2}{\|\mathbf{z}\|^2} + \frac{(\mathbf{a}_i^T \mathbf{x})^2}{\|\mathbf{x}\|^2} - 2 \frac{\mathbf{a}_i^T \mathbf{x} \mathbf{a}_i^T \mathbf{z}}{\|\mathbf{x}\| \|\mathbf{z}\|} \cos \theta)$, we use a simple formula $(|\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}|)^2$ to estimate this quantity.

Except the update rule, the form of the tanh weighted wirtinger flow also inspires us to propose a new initialization algorithm. Suppose $\mathbf{z} = \mathbf{x}$, we expect $\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^T \mathbf{x} \tanh \frac{y_i}{\sigma^2} \approx \mathbf{x}$, namely, \mathbf{x} is the leading eigenvector of the matrix $\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^T \tanh \frac{y_i}{\sigma^2}$. In the initialization step, we can replace the error term σ with a crude estimation. In conclusion, our new phasing algorithm is comprised of a novel initialization algorithm, a new form of wirtinger flow and a nesterov accelerated gradient update rule [13, 1]. We summarize all those components in algorithm 1. Moreover, our new formulation can also be used to interpret the

Algorithm 1: Tanh wirtinger Flow

Input : Measurements $\{y_i | 1 \leq i \leq m\}$ and sampling vectors $\{\mathbf{a}_i | 1 \leq i \leq m\}$; Initialization scale factor α , momentum μ and step size s .

Initialization: Drawn \mathbf{z}_0^0 from $\mathcal{N}(\mathbf{0}, \mathbf{I})$, and normalize it as $\mathbf{z}_0^0 = \frac{\mathbf{z}_0^0}{\|\mathbf{z}_0^0\|}$. Set

$$\hat{y} = \frac{1}{m} \sum_{i=1}^m y_i.$$

for $t = 1 : T_i$ **do**

$$\mathbf{z}_0^t = \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}_0^{t-1} \tanh \frac{y_i}{\alpha \hat{y}}$$

$$\mathbf{z}_0^t = \frac{\mathbf{z}_0^t}{\|\mathbf{z}_0^t\|}$$

end

Set $\mathbf{z}_0 = \mathbf{z}_0^{T_i}$.

Refinement : Set $\mathbf{v}_0 = \mathbf{0}$.

for $t = 1 : T_r$ **do**

$$\hat{\sigma}_i^2 = (|\mathbf{a}_i^T \mathbf{z}_{t-1}| - \sqrt{y_i})^2$$

$$\nabla l_t = \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (\mathbf{a}_i^T \mathbf{z}_{t-1} - \sqrt{y_i} \tanh \frac{\mathbf{a}_i^T \mathbf{z}_{t-1} \sqrt{y_i}}{\hat{\sigma}_i^2})$$

$$\mathbf{v}_t = \mu \mathbf{v}_{t-1} - s \nabla l_t$$

$$\mathbf{z}_t = \mathbf{z}_{t-1} - \mu \mathbf{v}_{t-1} + (1 + \mu) * \mathbf{v}_t$$

end

Output : \mathbf{z}_{T_r}

working mechanism of truncated wirtinger flow method [5]. From equation 1, we have

$$\epsilon_i = \frac{y_i - f_i^2}{2f_i} - \frac{\epsilon_i^2}{2f_i}.$$

As long as the term $\epsilon_i^2 \ll y_i - f_i^2, \frac{y_i - f_i^2}{2f_i}$ can be seen as a good estimator for ϵ_i , using the fact that $\mathbb{E} \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^T = \mathbf{I}$, thus we have

$$\mathbb{E} \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \frac{y_i - f_i^2}{2f_i} \approx \mathbb{E} \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i \epsilon_i = \mathbf{h},$$

which is a reminiscent of the interpretation about the truncated wirtinger flow.

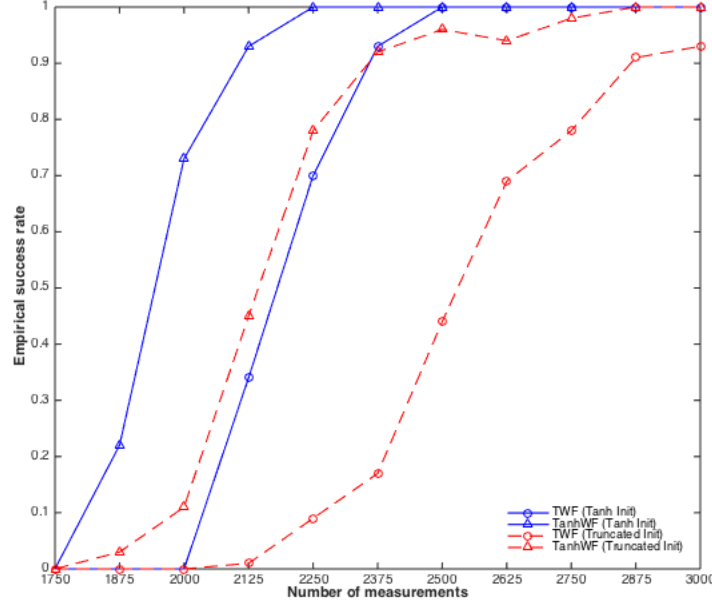
A final remark in this section is that though the tanh operator has been discovered in crystallographic community using likelihood heuristic for a long time, there is no theoretical justification for the effect of this operator. In this paper, we are aiming to provide a rigorous analysis for the tanh weighted wirtinger flow.

3. NUMERICAL EXPERIMENT

In this section, we report the numerical simulation result to demonstrate the effectiveness of our initialization method and update rule. In all the simulations performed in this paper, we used the following parameter settings: for the initialization stage, we used 100 power iterations; for truncated spectral initialization, we set the trimming threshold $\alpha_y = 3$; for tanh weighted spectral initialization, the scale factor α was set to be 4; for TWF, we adopted the default parameters used in [5] when calculating the gradient; we set the maximum number of updates to be 500, the learning rate to be 0.15 and the momentum to be 0.7. We compared different initialization methods by measuring the empirical success rates on random systems of quadratic equations with different data to parameter ratio. In these tests, we fixed the number of unknowns to be 1000 while varying the number of measurements from 1750 to 3000 with a step size of 125, and we considered a system as solved if the minimum relative error in the optimization process was smaller than 0.01. For each method and number of measurements, we conducted 100 trials in which systems to be solved were randomly generated each time. The TWF and TanhWF methods were compared from two aspects, which were the success rate and the convergence rate. When testing the TWF method, we replaced the gradient calculation formula in algorithm 1 with the formula defined in TWF. The final results are presented in figure 1 and 2. As it can be seen from figure 1, the tanh weighted spectral initialization significantly improved the success rate of solving quadratic systems. Besides, the TanhWF method also enjoys higher success rate comparing with TWF method when using the same initialization method. Figure 2 shows the ensembles of relative error histories for different methods on random quadratic systems with 1000 unknowns and 2500 measurements. It can be seen that the TanhWF method has superior stability and faster rate of convergence in solving these systems.

4. THEORETICAL ANALYSIS OF THE TANHWF METHOD

In this section, we will perform the convergence analysis for the TanhWF by verifying that it satisfies the regularity condition proposed in [3]. To establish the regularity condition, we need to bound two quantities, $-\langle \mathbf{h}, \frac{1}{2m} \nabla l \rangle$ and $\|\frac{1}{2m} \nabla l\|$, which are the curvature and smoothness of target function l , respectively. Given $-\langle \mathbf{h}, \frac{1}{2m} \nabla l \rangle \geq c \|\mathbf{h}\|^2$ and $\|\frac{1}{2m} \nabla l\|^2 \leq C \|\mathbf{h}\|^2$ hold with high probability, we have the following main theorem.

FIGURE 1. Empirical Success Rate for Quadratic System with 10^3 Unknowns

Theorem 1. *In the noiseless setting, there exist some universal constants $0 < \rho_0 < 1$ and c_1, c_2 such that with probability exceeding $1 - c_1 \exp[-c_2 m]$,*

$$\text{dist}^2(\mathbf{z} + \frac{\mu}{m} \nabla l(\mathbf{z}), \mathbf{x}) \leq (1 - \rho_0) \text{dist}^2(\mathbf{z}, \mathbf{x})$$

holds for all \mathbf{x}, \mathbf{h} satisfying $\|\mathbf{h}\| \leq \|\mathbf{x}\|$ and not in the set $\frac{\mathbf{x}^T \mathbf{h}}{\|\mathbf{x}\| \|\mathbf{h}\|} \in [0.6, 1] \cap \frac{\|\mathbf{h}\|}{\|\mathbf{x}\|} \in [0.6, 1]$ as long as μ is sufficiently small.

Proof. Without loss of generality, suppose $\|\mathbf{x} - \mathbf{z}\|^2 \leq \|\mathbf{x} + \mathbf{z}\|^2$, then we have

$$\begin{aligned} \text{dist}^2(\mathbf{z} + \frac{\mu}{m} \nabla l(\mathbf{z}), \mathbf{x}) &= \|\mathbf{h} - \frac{\mu}{m} \nabla l(\mathbf{z})\|^2 \\ &= \|\mathbf{h}\|^2 - 2\mu \langle \mathbf{h}, \frac{\nabla l(\mathbf{z})}{m} \rangle + \mu^2 \|\frac{\nabla l(\mathbf{z})}{m}\|^2 \\ &\leq \|\mathbf{h}\|^2 (1 - 4c\mu + 4C\mu^2). \end{aligned}$$

Therefore, as long as $1 - 4c\mu + 4C\mu^2 \leq 1$, namely, $0 \leq \mu \leq \frac{c}{C}$, the gradient update is contractive. In the following sections, we will prove that the curvature and smoothness conditions which lead to the last inequality hold with probability $1 - c \exp[-c' m]$ in the region stated in the theorem 1, thus completing the proof.

4.1. Proof of the local curvature condition. We first verify that the curvature satisfies the lower bound,

$$-\langle \mathbf{h}, \frac{1}{2m} \nabla^2 l \rangle \geq c \|\mathbf{h}\|^2,$$

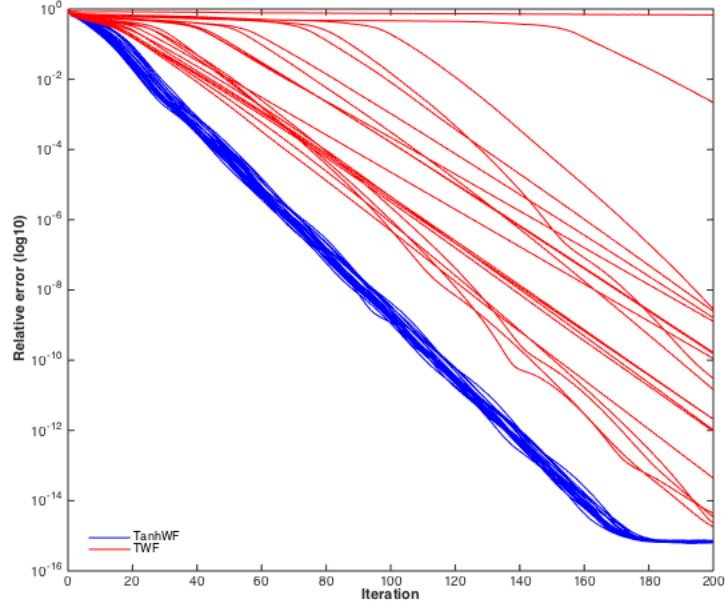


FIGURE 2. Relative Error v.s. Iteration for Different wirtinger Flow Methods

where c is a constant. We rewrite this quantity to strengthen its connection with $\|\mathbf{h}\|$,

$$\begin{aligned}
 -\langle \mathbf{h}, \frac{1}{2m} \nabla l \rangle &= \frac{1}{m} \sum_{i=1}^m \mathbf{h}^T \mathbf{a}_i [\mathbf{a}_i^T \mathbf{x} \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2} - \mathbf{a}_i^T \mathbf{z}] \\
 (4) \qquad \qquad \qquad &= \frac{1}{m} \sum_{i=1}^m \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{h} - (1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x}.
 \end{aligned}$$

The first term in equation 4 can be bounded using standard result since $\mathbf{a}_i^T \mathbf{h}$ is a simple gaussian random variable. It then boils down to showing that $\frac{1}{m} \sum_{i=1}^m (1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x} \leq \|\mathbf{h}\|^2$ holds with high probability. To investigate a random variable with complex form like $(1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x}$, we apply a dyadic decomposition to $\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}$ as in the proof of proposition 2.1.9 in [16] to simplify its structure, namely,

$$(5) \qquad \qquad \qquad \mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z} = \sum_{n=-\infty}^{\infty} (\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z})_n,$$

where $(\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z})_n = \mathbf{I}(a_n \sigma_i^2 \leq \mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z} \leq b_n \sigma_i^2) \mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}$, $a_n \sigma_i^2$ and $b_n \sigma_i^2$ specifies the interval of a dyadic component, and \mathbf{I} denoting the indicator function. Conditioned on one of those components, we can transform the bounds for $\frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{(|\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}|)^2}$

into linear inequalities between $\mathbf{a}_i^T \mathbf{x}$ and $\mathbf{a}_i^T \mathbf{h}$. Specifically, given

$$(6) \quad a \leq \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{(|\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}|)^2} \leq b,$$

letting $\mathbf{a}_i^T \mathbf{z} = \mathbf{a}_i^T \mathbf{x} - \mathbf{a}_i^T \mathbf{h}$, substitution into equation 6 yields

$$(7) \quad a \leq \frac{(\mathbf{a}_i^T \mathbf{x})^2 - \mathbf{a}_i^T \mathbf{h} \mathbf{a}_i^T \mathbf{x}}{(|\mathbf{a}_i^T \mathbf{x}| - |\mathbf{a}_i^T \mathbf{z}|)^2} \leq b.$$

For $0 < a < b$, $\mathbf{a}_i^T \mathbf{x}$ has the same sign as $\mathbf{a}_i^T \mathbf{z}$, thus the denominator in equation 7 equals with $(\mathbf{a}_i^T \mathbf{h})^2$. Assuming $\mathbf{a}_i^T \mathbf{h} > 0$, solving inequalities 7 gives

$$(8) \quad \frac{1 + \sqrt{4a+1}}{2} \mathbf{a}_i^T \mathbf{h} \leq \mathbf{a}_i^T \mathbf{x} \leq \frac{1 + \sqrt{4b+1}}{2} \mathbf{a}_i^T \mathbf{h},$$

$$(9) \quad \frac{1 - \sqrt{4b+1}}{2} \mathbf{a}_i^T \mathbf{h} \leq \mathbf{a}_i^T \mathbf{x} \leq \frac{1 - \sqrt{4a+1}}{2} \mathbf{a}_i^T \mathbf{h}.$$

We denote inequality 8 as event $\xi_{a,b}^1$ and inequality 9 as event $\xi_{a,b}^2$. For $a < b < 0$, $\mathbf{a}_i^T \mathbf{x}$ and $\mathbf{a}_i^T \mathbf{z}$ have opposite signs. Hence, the denominator in equation 7 can be expressed as $(2\mathbf{a}_i^T \mathbf{x} - \mathbf{a}_i^T \mathbf{h})^2$. Letting $\mathbf{a}_i^T \mathbf{h} > 0$, solving inequalities 7 gives

$$(10) \quad \frac{1 - 4a + \sqrt{1 - 4a}}{2(1 - 4a)} \mathbf{a}_i^T \mathbf{h} \leq \mathbf{a}_i^T \mathbf{x} \leq \frac{1 - 4b + \sqrt{1 - 4b}}{2(1 - 4b)} \mathbf{a}_i^T \mathbf{h},$$

$$(11) \quad \frac{1 - 4b - \sqrt{1 - 4b}}{2(1 - 4b)} \mathbf{a}_i^T \mathbf{h} \leq \mathbf{a}_i^T \mathbf{x} \leq \frac{1 - 4a - \sqrt{1 - 4a}}{2(1 - 4a)} \mathbf{a}_i^T \mathbf{h}.$$

Inequality 10 is denoted as event $\xi_{a,b}^3$, and inequality 11 is denoted as event $\xi_{a,b}^4$. If $\mathbf{a}_i^T \mathbf{h} \leq 0$, the left and right sides of the above inequalities should be switched. With these relations in hand, we can bound $(1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x} \mathbf{I}(a \leq \mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z} \leq b)$ using only $\mathbf{a}_i^T \mathbf{h}$ and a, b , thus simplifying the form of random variable to be investigated. We have the following proposition for calculating the probability of event $\xi_{a,b}^1 \cap \mathbf{a}_i^T \mathbf{h} \geq 0$. With a little abuse of notation, we will denote $\xi_{a,b}^1 \cap \mathbf{a}_i^T \mathbf{h} \in [x, y]$ as $\xi_{a,b}^1$ with the range of $\mathbf{a}_i^T \mathbf{h}$ to be inferred from the context.

Proposition 2. Suppose \mathbf{a}_i is a random Gaussian vector where each element has zero mean and unit variance, the probability of event $\xi_{a,b}^1$, given $b > a \geq 0$ and $\mathbf{a}_i^T \mathbf{h} \geq 0$, equals

$$(12) \quad P(\mathbf{I}(\xi_{a,b}^1) = 1) = \int_0^\infty \frac{e^{-\frac{t^2}{2\|\mathbf{h}\|^2}}}{2\pi\|\mathbf{h}\|} \int_{(\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}})^{-1}(\frac{-\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|^2} + b_+)^t}^{(\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}})^{-1}(\frac{-\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|^2} + a_+)^t} e^{-\frac{c^2}{2}} dc dt,$$

where $a_+ = \frac{1 + \sqrt{4a+1}}{2}$, and $b_+ = \frac{1 + \sqrt{4b+1}}{2}$.

The probabilities of other events can be expressed in the similar forms by replacing a_+ and b_+ with the corresponding values in those inequalities.

4.1.1. *Proof of the Proposition 2.* Proposition 2 can be derived as following. Set $\mathbf{a}_i^T \mathbf{h} = t$, denote a nonzero component of \mathbf{h} as h_n and the remaining component of \mathbf{h} as $\mathbf{h}_{\setminus n}$, and the corresponding components of \mathbf{a}_i as $\mathbf{a}_{i, \setminus n}$, we have $a_n = \frac{t - \mathbf{a}_{i, \setminus n}^T \mathbf{h}_{\setminus n}}{h_n}$. We can rewrite equation 8 as $(\mathbf{x}_{\setminus n}^T - \frac{x_n}{h_n} \mathbf{h}_{\setminus n}^T) \mathbf{a}_{i, \setminus n} < b_+ t - \frac{x_n t}{h_n}$ and

$(\mathbf{x}_{\setminus n}^T - \frac{x_n}{h_n} \mathbf{h}_{\setminus n}^T) \mathbf{a}_{i,\setminus n} > a_+ t - \frac{x_n t}{h_n}$. Denote $\mathbf{x}_{\setminus n} - \frac{x_n}{h_n} \mathbf{h}_{\setminus n}$ as \mathbf{w} and $b_+ t - \frac{x_n t}{h_n}$ as $f(b)$, and $a_+ t - \frac{x_n t}{h_n}$ as $f(a)$. We can then write the probability of event $\xi_{a,b}^1 \cap t > 0$ as

$$\begin{aligned} P(\mathbf{I}(\xi_{a,b}^1)) &= \int_0^\infty \int_{\mathbf{w}^T \mathbf{a}_{i,\setminus n} < f(b)} \frac{e^{-\frac{1}{2}(\frac{t - \mathbf{a}_{i,\setminus n}^T \mathbf{h}_{\setminus n}}{h_n})^2}}{h_n \sqrt{2\pi}^{\frac{n}{2}}} e^{-\frac{1}{2}\|\mathbf{a}_{i,\setminus n}\|^2} d\mathbf{a}_{i,\setminus n} dt \\ &= \int_0^\infty \int_{\mathbf{w}^T \mathbf{a}_{i,\setminus n} < f(b)} \frac{e^{-\frac{1}{2}(\frac{t^2}{h_n^2} + \frac{\|\mathbf{a}_{i,\setminus n}^T \mathbf{h}_{\setminus n}\|^2}{h_n^2} - \frac{2t\mathbf{a}_{i,\setminus n}^T \mathbf{h}_{\setminus n}}{h_n} + \|\mathbf{a}_{i,\setminus n}\|^2)}}{h_n \sqrt{2\pi}^{\frac{n}{2}}} d\mathbf{a}_{i,\setminus n} dt \\ &= \int_0^\infty \int_{\mathbf{w}^T \mathbf{a}_{i,\setminus n} < f(b)} \frac{e^{-\frac{t^2}{2\|\mathbf{h}\|^2}}}{h_n \sqrt{2\pi}^{\frac{n}{2}}} e^{-\frac{1}{2}(\mathbf{a}_{i,\setminus n} - \frac{t\mathbf{h}_{\setminus n}}{\|\mathbf{h}\|^2})^T (\mathbf{I} + \frac{\mathbf{h}_{\setminus n} \mathbf{h}_{\setminus n}^T}{h_n^2}) (\mathbf{a}_{i,\setminus n} - \frac{t\mathbf{h}_{\setminus n}}{\|\mathbf{h}\|^2})} d\mathbf{a}_{i,\setminus n} dt. \end{aligned}$$

By setting $\mathbf{a}_{i,\setminus n} - \frac{t\mathbf{h}_{\setminus n}}{\|\mathbf{h}\|^2}$ as $\mathbf{b}_{i,\setminus n}$, the boundary conditions of the above integral are transformed into $\mathbf{w}^T \mathbf{b}_{i,\setminus n} < -\frac{t\mathbf{h}_{\setminus n}^T \mathbf{w}}{\|\mathbf{h}\|^2} + f(b)$ and $\mathbf{w}^T \mathbf{b}_{i,\setminus n} > -\frac{t\mathbf{h}_{\setminus n}^T \mathbf{w}}{\|\mathbf{h}\|^2} + f(a)$, and the integral can be written as

$$P(\mathbf{I}(\xi_{a,b}^1)) = \int_0^\infty \int_{\mathbf{w}^T \mathbf{a}_{i,\setminus n} < f(b)} \frac{e^{-\frac{t^2}{2\|\mathbf{h}\|^2}}}{h_n \sqrt{2\pi}^{\frac{n}{2}}} e^{-\frac{1}{2}\mathbf{b}_{i,\setminus n}^T (\mathbf{I} + \frac{\mathbf{h}_{\setminus n} \mathbf{h}_{\setminus n}^T}{h_n^2}) \mathbf{b}_{i,\setminus n}} d\mathbf{b}_{i,\setminus n} dt.$$

Let $\mathbf{c}_{i,\setminus n} = (\mathbf{I} + \frac{\mathbf{h}_{\setminus n} \mathbf{h}_{\setminus n}^T}{h_n^2})^{\frac{1}{2}} \mathbf{b}_{i,\setminus n}$ and $\mathbf{w}' = (\mathbf{I} + \frac{\mathbf{h}_{\setminus n} \mathbf{h}_{\setminus n}^T}{h_n^2})^{-\frac{1}{2}} \mathbf{w}$. The determinant of matrix $\mathbf{I} + \frac{\mathbf{h}_{\setminus n} \mathbf{h}_{\setminus n}^T}{h_n^2}$ is $\frac{\|\mathbf{h}\|^2}{h_n^2}$, and its inverse is $\mathbf{I} - \frac{\mathbf{h}_{\setminus n} \mathbf{h}_{\setminus n}^T}{\|\mathbf{h}\|^2}$. The integral can be further simplified as

$$P(\mathbf{I}(\xi_{a,b}^1)) = \int_0^\infty \int_{\mathbf{w}'^T \mathbf{c}_{i,\setminus n} < -\frac{t\mathbf{h}_{\setminus n}^T \mathbf{w}}{\|\mathbf{h}\|^2} + f(b)} \frac{e^{-\frac{t^2}{2\|\mathbf{h}\|^2}}}{\sqrt{2\pi}^{\frac{n}{2}} \|\mathbf{h}\|} e^{-\frac{\mathbf{c}_{i,\setminus n}^T \mathbf{c}_{i,\setminus n}}{2}} d\mathbf{c}_{i,\setminus n} dt.$$

The last step is rotating the vector $\mathbf{c}_{i,\setminus n}$ with an orthogonal matrix \mathbf{U} whose first row is of the form $\frac{\mathbf{w}'}{\|\mathbf{w}'\|}$, which is the unit vector aligning with \mathbf{w}' . In addition, the norm of \mathbf{w}' is $\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}}$, and $-\frac{t\mathbf{h}_{\setminus n}^T \mathbf{w}}{\|\mathbf{h}\|^2} - \frac{x_n t}{h_n}$ can be reduced to $-\frac{\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|^2} t$. Combining these results and integrating out other components give the integral of the form shown in equation 12. Another way to prove this proposition is first deducing the joint probability of $\mathbf{a}_i^T \mathbf{h}$ and $\mathbf{a}_i^T \mathbf{x}$. Then, integrating over the region, $\mathbf{a}_i^T \mathbf{x} \in [a_+ \mathbf{a}_i^T \mathbf{h}, b_+ \mathbf{a}_i^T \mathbf{h}] \cap \mathbf{a}_i^T \mathbf{h} > 0$, leads to the probability of event $\xi_{a,b}^1$.

The correctness of proposition 12 can be easily verified as following. Suppose \mathbf{a}_i satisfies a linear inequality $\mathbf{a}_i^T \mathbf{x} \leq \mathbf{a}_i^T \mathbf{h} + b$, which is denoted as event ξ_b , applying proposition 12 leads to the probability of this event, that is,

$$(13) \quad P(\mathbf{I}(\xi_b) = 1) = \int_{-\infty}^\infty \frac{e^{-\frac{t^2}{2\|\mathbf{h}\|^2}}}{2\pi\|\mathbf{h}\|} \int_{-\infty}^{(\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}})^{-1}(\frac{-\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|^2} t + t + b)} e^{-\frac{c^2}{2}} dc dt.$$

We can also calculate the probability of event ξ_b by recognizing that the linear inequality restricts the eligible gaussian vectors lie below the hyperplane $\mathbf{a}_i^T (\mathbf{x} - \mathbf{h}) - b = 0$. According to [7], the probability of gaussian random variable in such region is

$$(14) \quad P(\mathbf{I}(\xi_b) = 1) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\frac{b}{\|\mathbf{x} - \mathbf{h}\|}} \exp(-\frac{x^2}{2}) dx.$$

We then need to show that equation 13 and equation 14 are equivalent. Let $x = \frac{t}{\|\mathbf{h}\|}$, we have

$$P(\mathbf{I}(\xi_b) = 1) = \int_{-\infty}^{\infty} \frac{e^{-\frac{x^2}{2}}}{2\pi} \int_{-\infty}^{\infty} (\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}})^{-1} \left(\frac{-\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|} x + \|\mathbf{h}\| x + b \right) e^{-\frac{c^2}{2}} dc dx.$$

Let $\phi(x) = \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}}$ and $\Phi(x) = \int_{-\infty}^x \phi(t) dt$, we can transform the above integral into

$$P(\mathbf{I}(\xi_b) = 1) = \int_{-\infty}^{\infty} \phi(x) \Phi(a + cx) dx,$$

where $a = \frac{b}{\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}}}$, and $c = \frac{-\frac{\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|} + \|\mathbf{h}\|}{\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}}}$. Using the integral identity $\int_{-\infty}^{\infty} \Phi(a + cx) \phi(x) = \Phi\left(\frac{a}{\sqrt{1+c^2}}\right)$ [14], we have

$$P(\mathbf{I}(\xi_b) = 1) = \Phi\left(\frac{b}{\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2} + (-\frac{\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|} + \|\mathbf{h}\|)^2}}\right) = \Phi\left(\frac{b}{\|\mathbf{x} - \mathbf{h}\|}\right),$$

thus showing the equivalence between equation 13 and equation 14.

4.1.2. *Proof of the $o(\|\mathbf{h}\|^2)$ expectation of $(1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x}$.* In this section, we denote the random variable $(1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x} \mathbf{I}(a < \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2} < b)$ as X_n , the factors $\frac{1+\sqrt{4a+1}}{2}$ and $\frac{1-4a+\sqrt{1-4a}}{2(1-4a)}$ as a_+ , and the factors $\frac{1-\sqrt{4a+1}}{2}$ and $\frac{1-4a-\sqrt{1-4a}}{2(1-4a)}$ as a_- . Other factors involving b are also named with the same convention. For $\frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2} > 0$, we decompose $\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}$ as

$$(15) \quad \mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z} \mathbf{I}\left(\frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2} > 0\right) = \sum_{n=-\infty}^{\infty} (\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z})_n,$$

where $(\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z})_n = \mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z} \mathbf{I}(2^{n-1} < \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2} < 2^n)$. Hence, we have $a = 2^{n-1}$ and $b = 2^n$, substitutions into equation 8 leads to

$$(16) \quad a_+ = \frac{1 + \sqrt{2^{n+1} + 1}}{2}, \quad b_+ = \frac{1 + \sqrt{2^{n+2} + 1}}{2}.$$

More suggestively, equation 16 can be summarized into a general formula, $f_n = \frac{1+\sqrt{2^{n+2}+1}}{2}$, which makes $a_+ = f_{n-1}$ and $b_+ = f_n$. On the event $\xi_{a,b}^1$ and $\mathbf{a}_i^T \mathbf{h} > 0$, since $f_n > 1$, we have $X_n^+ \leq (1 - \tanh 2^{n-1}) f_n (\mathbf{a}_i^T \mathbf{h})^2$, $\forall n \in [-\infty, \infty]$, where the plus superscript indicates X_n is conditioned on event $\xi_{a,b}^1$. Consequently, according to proposition 2, the expectation of X_n^+ is bounded by

$$(17) \quad \frac{\mathbb{E}(X_n^+)}{(1 - \tanh 2^{n-1}) f_n} \leq \int_0^{\infty} t^2 \frac{e^{-\frac{t^2}{2\|\mathbf{h}\|^2}}}{2\pi \|\mathbf{h}\|} \int_{(\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}})^{-1} (\frac{-\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|} + f_n)t}^{(\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}})^{-1} (\frac{-\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|} + f_{n-1})t} e^{-\frac{c^2}{2}} dc dt.$$

Letting $\sqrt{\|\mathbf{x}\|^2 - \frac{|\mathbf{x}^T \mathbf{h}|^2}{\|\mathbf{h}\|^2}} = w$, $\frac{\|\mathbf{h}\|}{w} (-\frac{\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|} + f_n) = u_n$ and $\frac{\|\mathbf{h}\|}{w} (-\frac{\mathbf{x}^T \mathbf{h}}{\|\mathbf{h}\|} + f_{n-1}) = u_{n-1}$. If the angle between \mathbf{x} and \mathbf{h} be θ , u_n can be reduced to

$$u_n = \frac{1}{\sqrt{1 - \cos^2 \theta}} (-\cos \theta + \frac{\|\mathbf{h}\|}{\|\mathbf{x}\|} f_n).$$

As the r.h.s of equation 17 after integration is $\frac{\|\mathbf{h}\|^2}{2\pi}(\frac{u_n}{1+u_n^2} - \frac{l_{n-1}}{1+l_{n-1}^2} + \arctan u_n - \arctan l_{n-1})$, we thus have the following upper bound for $\mathbb{E}(X_n^+)$,

$$(18) \quad \frac{\mathbb{E}(X_n^+)}{\|\mathbf{h}\|^2} \leq (1 - \tanh 2^{n-1}) \frac{f_n}{2\pi} \left(\frac{u_n}{1+u_n^2} - \frac{u_{n-1}}{1+u_{n-1}^2} + \arctan u_n - \arctan u_{n-1} \right).$$

For $n \geq 1$, the convergence rate of the above bound is at least of the order $e^{-2^n} 2^n$ since $\frac{u}{1+u^2} + \arctan u$ is a bounded function. If $n \leq 0$, $1 - \tanh 2^{n-1}$ will not decay exponentially. However, since $f_n \leq f_0$ for all $n \leq 0$, the expectation is upper bounded by

$$(19) \quad \frac{\mathbb{E}(X_n^+)}{\|\mathbf{h}\|^2} \leq \frac{f_0}{2\pi} \left(\frac{u_n}{1+u_n^2} - \frac{u_{n-1}}{1+u_{n-1}^2} + \arctan u_n - \arctan u_{n-1} \right).$$

Summing up equation 19 for all $n \leq 0$ gives rise to

$$(20) \quad \frac{\mathbb{E}(X_{-\infty,0}^+)}{\|\mathbf{h}\|^2} \leq \frac{f_0}{2\pi} \left(\frac{u_0}{1+u_0^2} - \frac{u_{-\infty}}{1+u_{-\infty}^2} + \arctan u_0 - \arctan u_{-\infty} \right),$$

which is a function with bounded values. We plotted the ratio between the upper bound in equation 20 and $\frac{\|\mathbf{h}\|}{\|\mathbf{x}\|}$ on a grid with $\frac{\|\mathbf{h}\|}{\|\mathbf{x}\|} \in [0, 1]$ and $\cos \theta \in [-1, 1]$. As it's shown in figure 3, $\frac{\mathbb{E}(X_{-\infty,0}^+)}{\|\mathbf{h}\|^2} \leq \frac{\|\mathbf{h}\|}{\|\mathbf{x}\|}$ holds in most regions.

On the event $\xi_{a,b}^2$ and $\mathbf{a}_i^T \mathbf{h} > 0$, the general formula for a_- and b_- can be expressed as $f_n = \frac{1-\sqrt{2^{n+2}+1}}{2}$ with $a_- = f_{n-1}$ and $b_- = f_n$. Since $f_n < 0$, the expectation is upper bounded by

$$(21) \quad \frac{\mathbb{E}(X_n^-)}{\|\mathbf{h}\|^2} \leq (1 - \tanh 2^n) \frac{f_{n-1}}{2\pi} \left(\frac{u_{n-1}}{1+u_{n-1}^2} - \frac{u_n}{1+u_n^2} + \arctan u_{n-1} - \arctan u_n \right).$$

We proceed to obtain the bounds for expectations conditioned on event $\xi_{a,b}^3$ or $\xi_{a,b}^4$. Conditioned on event $\xi_{a,b}^3$ and $\mathbf{a}_i^T \mathbf{h} > 0$, we employ a similar decomposition as equation 15 for $\frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2} < 0$, thus leading to $a = -2^n$, and $b = -2^{n-1}$. The general formula for a_+ and b_+ is $f_n = \frac{1}{2} + \frac{1}{2\sqrt{1+2^{n+1}}}$. In addition, we have $a_+ = f_{n+1}$ and $b_+ = f_n$. Since $f_n > 0$ on event $\xi_{a,b}^3$ and $\mathbf{a}_i^T \mathbf{h} > 0$, the upper bound of the expectation is similar to the bound in equation 18 and can be written as

$$\frac{\mathbb{E}(X_n^+)}{\|\mathbf{h}\|^2} \leq (1 + \tanh 2^n) \frac{f_n}{2\pi} \left(\frac{u_n}{1+u_n^2} - \frac{u_{n+1}}{1+u_{n+1}^2} + \arctan u_n - \arctan u_{n+1} \right),$$

On event $\xi_{a,b}^4$ and $\mathbf{a}_i^T \mathbf{h} > 0$, we have the general formula $f_n = \frac{1}{2} - \frac{1}{2\sqrt{1+2^{n+1}}} > 0$ with $a_- = f_{n+1}$ and $b_- = f_n$. thus resulting in the following upper bound,

$$\frac{\mathbb{E}(X_n^-)}{\|\mathbf{h}\|^2} \leq (1 + \tanh 2^n) \frac{f_{n+1}}{2\pi} \left(\frac{u_{n+1}}{1+u_{n+1}^2} - \frac{u_n}{1+u_n^2} + \arctan u_{n+1} - \arctan u_n \right).$$

The convergence behavior of the sum of those bounds can be analyzed using the same method as bounding the expectation on event $\xi_{a,b}^1$. To obtain a detailed view about how the size of bound changes with respect to $\frac{\|\mathbf{h}\|}{\|\mathbf{x}\|}$ and $\cos \theta$, we evaluated the sum of the upper bounds for all events from $n = -20$ to 20 over a grid with $\frac{\|\mathbf{h}\|}{\|\mathbf{x}\|} \in [0.01, 1]$ and $\cos \theta \in [-0.999, 0.999]$. Note that the upper bounds for the expectations conditioned on $\mathbf{a}_i^T \mathbf{h} < 0$ is the same as the upper bounds conditioned on $\mathbf{a}_i^T \mathbf{h} > 0$, we multiplied the sum by a factor of 2. The contour plot for the

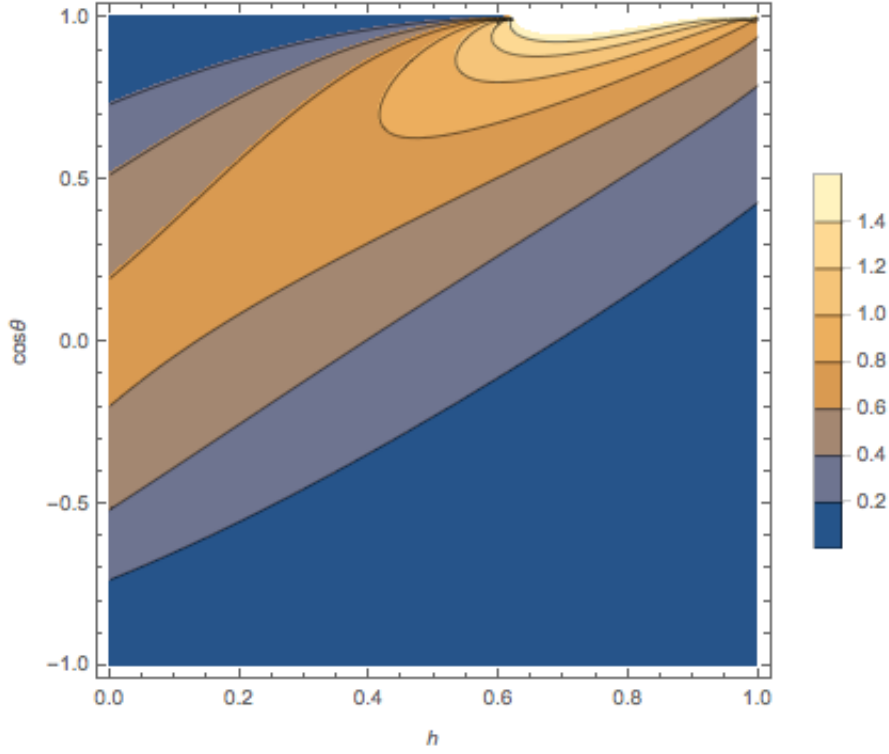


FIGURE 3. Contour plot for the upper bound of $\frac{\mathbb{E}(X_{-\infty,0}^+) \|\mathbf{x}\|^2}{\|\mathbf{h}\|^3}$ w.r.t different $\frac{\|\mathbf{h}\|}{\|\mathbf{x}\|}$ and $\cos \theta$

final sum minus 1 is shown in figure 4. It can be seen from figure 4 that the upper bound of $\frac{\mathbb{E}(X)}{\|\mathbf{h}\|^2} - 1$ decreases as the relative error becomes smaller and $\frac{\mathbb{E}(X)}{\|\mathbf{h}\|^2} - 1 < 0$ in most regions. Moreover, the upper bound for $\frac{\mathbb{E}(X) \|\mathbf{x}\|}{\|\mathbf{h}\|^3}$ is plotted in figure 5. Therefore, we conclude that $\mathbb{E}((1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x}) \leq 1.2 \frac{\|\mathbf{h}\|^3}{\|\mathbf{x}\|}$ holds for all $\|\mathbf{h}\|$ and $\cos \theta$ not in the set $\cos \theta \in [0.8, 1] \cap \frac{\|\mathbf{h}\|}{\|\mathbf{x}\|} \in [0.6, 1]$. Besides, we have $\mathbb{E}((1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x}) \leq 0.8 \|\mathbf{h}\|^2$ for all $\|\mathbf{h}\|, \cos \theta$ not in the set $\cos \theta \in [0.6, 1] \cap \frac{\|\mathbf{h}\|}{\|\mathbf{x}\|} \in [0.6, 1]$.

4.1.3. Proof of the sub-exponential concentration property of $(1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x}$.

The next step is to demonstrate that the random variable $(1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x}$ is tightly concentrated around its expectation. Based on the analysis in section 4.1.2, X_n can be bounded by $c_n (\mathbf{a}_i^T \mathbf{h})^2$, where c_n is a constant depends on events and the dyadic interval. We can obtain a universal bound for X_n by finding the supreme of the absolute value of c_n . On the event $\xi_{a,b}^1$, X_n is positive and upper bounded by $\frac{1 + \sqrt{2^{n+2} + 1}}{2} (1 - \tanh 2^{n-1})$; on the event $\xi_{a,b}^2$, X_n is negative and lower bounded by $\frac{1 - \sqrt{2^{n+2} + 1}}{2} (1 - \tanh 2^{n-1})$; on the event $\xi_{a,b}^3$, X_n is positive and upper

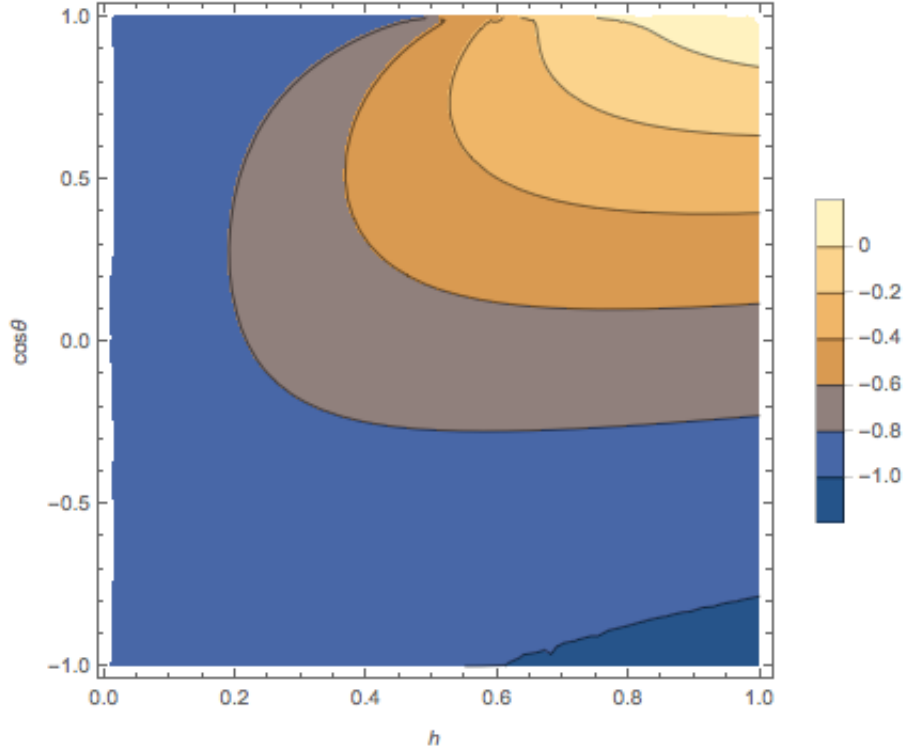


FIGURE 4. Contour plot for the upper bound of $\frac{\mathbb{E}(X)}{\|\mathbf{h}\|^2} - 1$ w.r.t different $\frac{\|\mathbf{h}\|}{\|\mathbf{x}\|}$ and $\cos \theta$

bounded by $(\frac{1}{2} + \frac{1}{2\sqrt{1+2^{n+1}}})(1 + \tanh 2^n)$; on the event $\xi_{a,b}^4$, X_n is positive and upper bounded by $(\frac{1}{2} - \frac{1}{2\sqrt{1+2^{n+2}}})(1 + \tanh 2^n)$. Combining all these bounds, we have

$$\sup_n |c_n| = (\frac{1}{2} + \frac{1}{2\sqrt{5}})(1 + \tanh 2) \approx 1.42,$$

thus leading to $|(1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x}| \leq 1.42 (\mathbf{a}_i^T \mathbf{h})^2$ for all \mathbf{a}_i . Since the sub-gaussian norm of $\mathbf{a}_i^T \mathbf{h}$ is $\sqrt{\frac{8}{3}} \|\mathbf{h}\|$, using lemma 2.7.6 in [17] results in the sub-exponential norm of the upper bound, namely, $\|1.42 (\mathbf{a}_i^T \mathbf{h})^2\|_{\psi_1} = 3.79 \|\mathbf{h}\|^2$. Denote $(1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2}) \mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x}$ as X_i , we have $\mathbb{E}(\exp(|X_i|/C)) \leq \mathbb{E}(\exp(1.42 (\mathbf{a}_i^T \mathbf{h})^2/C))$ for all $C > 0$. Hence, $\|X_i\|_{\psi_1} \leq 1.42 \|(\mathbf{a}_i \mathbf{h})^2\|_{\psi_1}$. Combining all these results, we conclude that the sub-exponential norm of $X_i - \mathbb{E}(X_i)$ is smaller than $C \|\mathbf{h}\|^2$. Using Bernstein's inequality [17], for every $t > 0$, we have

$$\mathbb{P}(|\frac{1}{m} \sum_{i=1}^m X_i - \mathbb{E}(X_i)| \geq t) \leq 2 \exp(-c \min(\frac{t}{C \|\mathbf{h}\|^2}, \frac{t^2}{C \|\mathbf{h}\|^4})),$$

where $c > 0, C > 0$ are absolute constants.

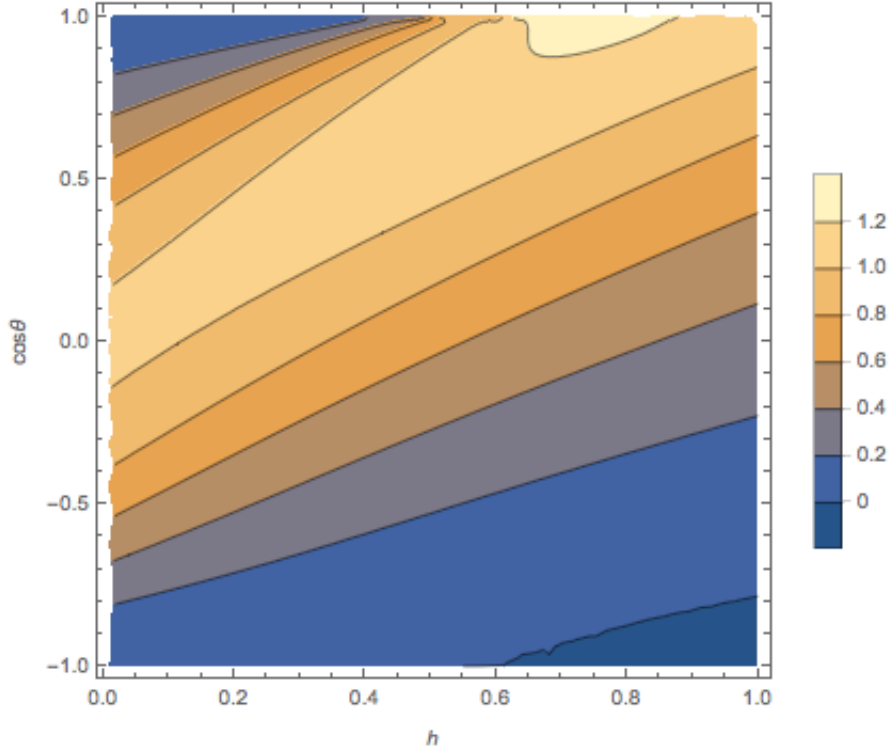


FIGURE 5. Contour plot for the upper bound of $\frac{\mathbb{E}(X)\|\mathbf{x}\|}{\|\mathbf{h}\|^3}$ w.r.t different $\frac{\|\mathbf{h}\|}{\|\mathbf{x}\|}$ and $\cos \theta$

4.2. Proof of the local smoothness condition. To fulfill establishing the regularity condition, we remains to verify that

$$\left\| \frac{1}{2m} \nabla l \right\|^2 = \frac{1}{m^2} \mathbf{v}^T \mathbf{M} \mathbf{v} \leq C \|\mathbf{h}\|^2,$$

where C is an absolute constant, $\mathbf{M} = \mathbf{A}^T \mathbf{A}$, $\mathbf{A} \equiv [\mathbf{a}_1, \dots, \mathbf{a}_n]$, and $\mathbf{v}^T \equiv [\mathbf{a}_1^T(\mathbf{z} - \mathbf{x} \tanh(\frac{\mathbf{x}^T \mathbf{a}_1 \mathbf{a}_1^T \mathbf{z}}{\sigma_1^2})), \dots, \mathbf{a}_n^T(\mathbf{z} - \mathbf{x} \tanh(\frac{\mathbf{x}^T \mathbf{a}_n \mathbf{a}_n^T \mathbf{z}}{\sigma_n^2}))]$, holds with high probability. An application of the inequality $\left\| \frac{1}{2m} \mathbf{A}^T \mathbf{v} \right\| \leq \frac{1}{2m} \|\mathbf{A}\| \|\mathbf{v}\|$ simplifies the terms to be considered for bounding the norm of gradients. Since we have $\|\mathbf{A}\| \leq \sqrt{m}(1 + \delta)$ from standard random matrix result, the problem then boils down to controlling $\frac{1}{\sqrt{m}} \|\mathbf{v}\|$. This in turns drives us to investigate the concentration property and expectation of $\mathbf{a}_i^T(\mathbf{z} - \mathbf{x} \tanh(\frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma_i^2}))$. We first rewrite it to gain some insights

about the connection between $\|\mathbf{v}\|^2$ and $\|\mathbf{h}\|^2$,

$$\begin{aligned}
 |\mathbf{a}_i^T(\mathbf{z} - \mathbf{x} \tanh(\frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}))|^2 &= (-\mathbf{a}_i^T \mathbf{h} + \mathbf{a}_i^T \mathbf{x} (1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}))^2 \\
 &= (\mathbf{a}_i^T \mathbf{h})^2 - 2\mathbf{a}_i^T \mathbf{h} \mathbf{a}_i^T \mathbf{x} (1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}) \\
 &\quad + (\mathbf{a}_i^T \mathbf{x})^2 (1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2})^2.
 \end{aligned}
 \tag{22}$$

In equation 22, the first term $(\mathbf{a}_i^T \mathbf{h})^2$ is a random variable with known property, and the second term $\mathbf{a}_i^T \mathbf{h} \mathbf{a}_i^T \mathbf{x} (1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2})$ has been investigated in previous section and has a sub-exponential norm of size $O(\sigma^2)$. We then proceed to show that the remaining term $(\mathbf{a}_i^T \mathbf{x})^2 (1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2})^2$ is also a random variable with $O(\sigma^2)$ sub-exponential norm. We use the same method as bounding the sub-exponential norm of $\mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x} (1 - \tanh(\frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}))$. To show that $\mathbb{E}((\mathbf{a}_i^T \mathbf{x})^2 (1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2})^2) \leq C\|\mathbf{h}\|^2$, it's enough to bound its sub-exponential norm since the expectation of a random variable is smaller than its sub-exponential norm up to an absolute constant [17]. Of course, we can use the method for bounding the expectation of $\mathbf{h}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{x} (1 - \tanh(\frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}))$ to obtain more precise bound. From section 4.1.3, we have $|(1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2}) \mathbf{a}_i^T \mathbf{x}| \leq 1.42|\mathbf{a}_i^T \mathbf{h}|$ for all \mathbf{a}_i . Denote $(\mathbf{a}_i^T \mathbf{x})^2 (1 - \tanh \frac{\mathbf{x}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z}}{\sigma^2})^2$ as X_i , then $X_i \leq 2.02(\mathbf{a}_i^T \mathbf{h})^2$ holds for all \mathbf{a}_i . We thus jump to the conclusion that X_i is a random variable with $O(\sigma^2)$ sub-exponential norm. We then confirm that $\frac{1}{2m}\|\nabla l\|^2 \leq C\|\mathbf{h}\|^2$ holds with Bernstein type tail bounds.

4.3. Initialization via tanh weighted spectral method. This section provides a theoretical analysis for the success of tanh weighted spectral initialization method. The leading eigenvector of the tanh weighted design matrix can be defined as

$$\mathbf{z} = \sup_{\mathbf{z} \in S^{n-1}} \frac{1}{m} \sum_{i=1}^m \mathbf{z}^T \mathbf{a}_i \mathbf{a}_i^T \mathbf{z} \tanh \frac{|\mathbf{a}_i^T \mathbf{x}|^2}{\alpha},$$

where S^{n-1} represents the unit sphere in \mathbb{R}^n . To show that the leading eigenvector \mathbf{z} is close to the true signal \mathbf{x} , we should prove that $\frac{1}{m} \sum_{i=1}^m (|\mathbf{a}_i^T \mathbf{x}|^2 - |\mathbf{a}_i^T \mathbf{z}|^2) \tanh \frac{|\mathbf{a}_i^T \mathbf{x}|^2}{\alpha} \geq 0$ holds with high probability when $\text{dist}(\mathbf{x} - \mathbf{z}) \geq \rho$, thus contradicting the assumption that \mathbf{z} maximizes the target function $\frac{1}{m} \sum_{i=1}^m |\mathbf{a}_i^T \mathbf{z}|^2 \tanh \frac{|\mathbf{a}_i^T \mathbf{x}|^2}{\alpha}$. Without loss of generality, we assume $\mathbf{x}, \mathbf{z} \in S^{n-1}$ since we can always absorbing the norm of \mathbf{x} into α by setting it to be $\frac{\alpha}{\|\mathbf{x}\|^2}$. This problem can be greatly simplified by leveraging its intrinsic rotation invariance. $\forall \mathbf{x}, \mathbf{z} \in S^{n-1}$, we can rotate and project them with a rotation projection matrix \mathbf{U} whose first row is \mathbf{x} , and second row is on the hyperplane spanned by \mathbf{x} and \mathbf{z} . Specifically, the matrix \mathbf{U} can be written as

$$\mathbf{U} = \begin{bmatrix} \mathbf{x}/\|\mathbf{x}\| \\ (\mathbf{z} - \frac{\mathbf{z}^T \mathbf{x}}{\|\mathbf{x}\|^2} \mathbf{x})/\|\mathbf{z} - \frac{\mathbf{z}^T \mathbf{x}}{\|\mathbf{x}\|^2} \mathbf{x}\| \end{bmatrix}.$$

We thus reduce the problem to \mathbb{R}^2 . Applying the transform leads to $\mathbf{x}' = \mathbf{U}\mathbf{x}$ which is the vector $[1, 0]$, and $\mathbf{z}' = \mathbf{U}\mathbf{z}$ which is the vector $[\cos \theta, \sin \theta]$ (θ is the angle between \mathbf{x} and \mathbf{z}). The inner product $\mathbf{a}_i^T \mathbf{x}$ is transformed to $\mathbf{a}_i^T \mathbf{U}^T \mathbf{x}'$, thus

prompting us to study the random vector $\mathbf{U}\mathbf{a}_i$. The distribution of gaussian random vector is invariant under rotation, while projecting the gaussian random vector from \mathbb{R}^n to \mathbb{R}^2 yields a gaussian random vector in \mathbb{R}^2 . The 2D gaussian random vector is of the form $r[\cos \phi, \sin \phi]$, where ϕ is uniformly distributed in $[0, 2\pi]$, and r is distributed according to $r \exp(-\frac{r^2}{2})$ in $[0, \infty]$. We can then rewrite the inner products as

$$\mathbf{a}_i^T \mathbf{x} = r \cos \phi, \mathbf{a}_i^T \mathbf{z} = r \cos(\phi - \theta).$$

Based on the above equations, we have

$$(23) \quad (|\mathbf{a}_i^T \mathbf{x}|^2 - |\mathbf{a}_i^T \mathbf{z}|^2) \tanh \frac{|\mathbf{a}_i^T \mathbf{x}|^2}{\alpha} = r^2(\cos^2 \phi - \cos^2(\phi - \theta)) \tanh \frac{r^2 \cos^2 \phi}{\alpha}.$$

We first obtain the expectation of 23. However, there is no analytical form for the integration of the above formula. We consider approximating $\tanh(x)$ with a new function $f(x) = 1 - \exp(-x)$ here. Thus we have,

$$\begin{aligned} \mathbb{E}((|\mathbf{a}_i^T \mathbf{x}|^2 - |\mathbf{a}_i^T \mathbf{z}|^2) f(\frac{|\mathbf{a}_i^T \mathbf{x}|^2}{\alpha})) &= \int_0^\infty \int_0^{2\pi} \frac{r^3}{2\pi} e^{-\frac{r^2}{2}} (\cos^2 \phi - \cos^2(\phi - \theta)) \\ &\quad (1 - \exp(\frac{-r^2 \cos^2 \phi}{\alpha})) dr d\phi \\ &= \int_0^\infty \exp(-(\frac{1}{\alpha} + 1)\frac{r^2}{2}) I_1(\frac{r^2}{2\alpha}) \sin^2 \theta r^3 dr \\ &= 2\sqrt{\frac{\alpha}{(2+\alpha)^3}} \sin^2 \theta, \end{aligned}$$

where $I_1(z)$ gives the modified Bessel function of the first kind. It's then easy to see that the expectation of $f(\frac{|\mathbf{a}_i^T \mathbf{x}|^2}{\alpha})$ weighted difference $|\mathbf{a}_i^T \mathbf{x}|^2 - |\mathbf{a}_i^T \mathbf{z}|^2$ has positive expectation. Moreover, for larger θ , that is, the estimation \mathbf{z} is far away from the signal \mathbf{x} , the expectation is bigger. To understand the concentration behavior of the random variable $r^2(\cos^2 \phi - \cos^2(\phi - \theta)) f(\frac{r^2 \cos^2 \phi}{\alpha})$, we turn to bound its sub-exponential norm, which can be written as

$$\begin{aligned} \|r^2(\cos^2 \phi - \cos^2(\phi - \theta)) f(\frac{r^2 \cos^2 \phi}{\alpha})\|_{\psi_1} &= \int_0^\infty \int_0^{2\pi} \exp(\frac{1}{\lambda} r^2 |\cos^2 \phi - \cos^2(\phi - \theta)|) \\ &\quad f(\frac{r^2 \cos^2 \phi}{\alpha}) \frac{r}{2\pi} \exp(-\frac{r^2}{2}) dr d\phi. \end{aligned}$$

It's obvious that this random variable has a $O(1)$ sub-exponential norm since $|(\cos^2 \phi - \cos^2(\phi - \theta)) f(\frac{r^2 \cos^2 \phi}{\alpha})| \leq 1$. Denote $(|\mathbf{a}_i^T \mathbf{x}|^2 - |\mathbf{a}_i^T \mathbf{z}|^2) f(\frac{|\mathbf{a}_i^T \mathbf{x}|^2}{\alpha})$ as X_i , we have the following tail bounds for all t ,

$$\mathbb{P}(|\frac{1}{m} \sum_{i=1}^m X_i - \mathbb{E}X_i| \geq t) \leq 2 \exp(-cmt).$$

5. CONCLUDING REMARKS

In this paper, we presented a new phase retrieval algorithm which employs the tanh activation function to weight the current estimation about the phase for each measurement. We have shown that the TanhWF method has higher success rate in solving random systems of quadratic equations than the TWF method when using the same initialization method and parameter update rule. In addition, we also proposed a new tanh weighted spectral initialization method which significantly

improved the success rate comparing with the truncated initialization method. We have proved that the TanhWF method satisfies the regularity condition for gaussian design matrix [3]. One possible future work is to extend our approach to complex-valued signals. As it shown by the Grothendieck's identity, the phase retrieval problem in \mathbb{R}^n space can be casted to \mathbb{R}^2 plane. Similarly, we might be able to investigate the complex-valued phasing problem in \mathbb{C}^n by casting it to \mathbb{C}^2 space.

REFERENCES

- [1] Yoshua Bengio, Nicolas Boulanger-Lewandowski, and Razvan Pascanu. Advances in optimizing recurrent networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 8624–8628. IEEE, 2013.
- [2] G Bricogne and JJ Irwin. Maximum-likelihood refinement of incomplete models with buster+tnt. *Macromolecular refinement. Proceedings of the CCP4 Study Weekend at Chester College*, pages 85–92, 1996.
- [3] Emmanuel J Candes, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via Wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, 2015.
- [4] Emmanuel J Candes, Thomas Strohmer, and Vladislav Voroninski. Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 66(8):1241–1274, 2013.
- [5] Yuxin Chen and Emmanuel Candes. Solving random quadratic systems of equations is nearly as easy as solving linear systems. In *Advances in Neural Information Processing Systems*, pages 739–747, 2015.
- [6] James R Fienup. Phase retrieval algorithms: a comparison. *Applied optics*, 21(15):2758–2769, 1982.
- [7] Daniel Fischer. <https://math.stackexchange.com/questions/556977/gaussian-integrals-over-a-half-space>. 2013.
- [8] Ralph W Gerchberg. A practical algorithm for the determination of the phase from image and diffraction plane pictures. *Optik*, 35:237–246, 1972.
- [9] Herbert A Hauptman. The phase problem of x-ray crystallography. *Reports on Progress in Physics*, 54(11):1427, 1991.
- [10] Ritesh Kolte and Ayfer zgr. Phase retrieval via incremental truncated Wirtinger flow. *arXiv preprint arXiv:1606.03196*, 2016.
- [11] VY Lunin, PV Afonine, and AG Urzhumtsev. Likelihood-based refinement. I. Irremovable model errors. *Acta Crystallographica Section A: Foundations of Crystallography*, 58(3):270–282, 2002.
- [12] Garib N Murshudov, Alexei A Vagin, and Eleanor J Dodson. Refinement of macromolecular structures by the maximum-likelihood method. *Acta Crystallographica Section D: Biological Crystallography*, 53(3):240–255, 1997.
- [13] Yurii Nesterov. A method of solving a convex programming problem with convergence rate $O(1/k^2)$. In *Soviet Mathematics Doklady*, volume 27, pages 372–376, 1983.
- [14] DB Owen. A table of normal integrals: A table. *Communications in Statistics-Simulation and Computation*, 9(4):389–419, 1980.
- [15] Navraj S Pannu and Randy J Read. Improved structure refinement through maximum likelihood. *Acta Crystallographica Section A: Foundations of Crystallography*, 52(5):659–668, 1996.
- [16] Terence Tao. *Topics in random matrix theory*, volume 132. American Mathematical Society Providence, RI, 2012.
- [17] Roman Vershynin. *High Dimensional Probability*. 2016.
- [18] Gang Wang, Georgios B Giannakis, and Yonina C Eldar. Solving systems of random quadratic equations via truncated amplitude flow. *arXiv preprint arXiv:1605.08285*, 2016.
- [19] Huishuai Zhang, Yuejie Chi, and Yingbin Liang. Provable Non-convex Phase Retrieval with Outliers: Median Truncated Wirtinger Flow. In *International conference on machine learning*, pages 1022–1031, 2016.