

# Robust Estimation for Class Averaging in Cryo-EM Single Particle Reconstruction

Chenxi Huang<sup>1</sup> and Hemant D. Tagare<sup>2</sup>

**Abstract**—Single Particle Reconstruction (SPR) for Cryogenic Electron Microscopy (cryo-EM) aligns and averages the images extracted from micrographs to improve the Signal-to-Noise ratio (SNR). Outliers compromise the fidelity of the averaging. We propose a robust cross-correlation-like w-estimator for combating the effect of outliers on the average images in cryo-EM. The estimator accounts for the natural variation of signal contrast among the images and eliminates the need for a threshold for outlier rejection. We show that the influence function of our estimator is asymptotically bounded. Evaluations of the estimator on simulated and real cryo-EM images show good performance in the presence of outliers.

## I. INTRODUCTION

Cryogenic Electron Microscopy (cryo-EM) is a method for reconstructing the three dimensional (3D) structure of macromolecular assemblies, referred to as particles. Single Particle Reconstruction (SPR) methods determine the 3D structure of a particle from its two dimensional (2D) projection images by iteratively 1) projecting an estimate of the 3D structure from different directions to obtain a set of 2D templates, 2) aligning images to the most similar template, 3) calculating an average image from the aligned set, and 4) updating the 3D structure from the average images by reconstruction algorithms.

The averaging step is critical to SPR; it is the step responsible for enhancing the Signal-to-Noise ratio (SNR) and the resolution of the reconstruction. However, the projection images obtained from electron micrographs during particle picking often contain a large number of outliers, which are due to particle fragments, contaminants, or pure noise [1]. Outliers also arise during alignment when images are aligned to incorrect templates. Outliers that survive particle picking and alignment propagate through averaging to the reconstruction algorithm and compromise the accuracy of the reconstruction. We propose a robust averaging method for cryo-EM to limit the effect of outliers.

Outlier rejection during particle picking is often realized by thresholding the cross-correlation between the images extracted from micrographs and one or more 2D templates (e.g., averages of a few manually selected images) [2]. Images with correlation coefficient less than the threshold are rejected as outliers. Although correlation has proven useful for outlier detection in cryo-EM, no theoretical justification

for the threshold or asymptotic robustness of these outlier detection methods is available. This paper addresses these shortcomings. Drawing on classical robust estimation theory [3], we propose a cross-correlation-like “w-estimator” for cryo-EM that calculates the average image by a weighted sum of the observations. The weight function of this w-estimator is an adaptation of cross-correlation to satisfy the requirement of asymptotic robustness.

From a statistical point of view, the average image is an estimate of the mean of an underlying distribution. Robust estimation also aims to estimate the mean when the sample contains outliers. The robustness of an estimator is quantified by the influence function (IF). The IF of an estimator at a point  $\mathbf{x}$  measures the effect of an outlier at  $\mathbf{x}$  on the estimate [4]. A robust estimator has a bounded IF at those  $\mathbf{x}$  where outliers are expected. Experience with real cryo-EM data suggests that outliers in cryo-EM images have the following characteristics: 1) a finite component along the average image, but 2) a low correlation with the average. We focus on the boundedness of the influence function for such outlier images. We show in this paper that our estimator is Fisher consistent, has a bounded influence function, and is robust to image contrast.

The rest of the paper is organized as follows: Section II contains the image model, the proposed estimator, the derivation of the influence function, and the analysis of robustness based on the influence function. Section III shows the results of our estimator for simulated and real cryo-EM data and Section IV is the conclusion. All the proofs are given in the appendix.

## II. METHOD

### A. Image Model

Assuming independent Gaussian white noise and taking into account image contrast variation, we define the model of an inlier image  $\mathbf{x}$  as:

$$\mathbf{x} = s\boldsymbol{\theta} + \mathbf{n} \quad (1)$$

where  $\boldsymbol{\theta} \in \mathbb{R}^p$  is the projected particle signal,  $p$  is the number of pixels in the image,  $\mathbf{n} \in \mathbb{R}^p$  is the Gaussian white noise with zero mean and standard deviation  $\sigma$  and  $s$  is the scale factor modeling image contrast;  $s$  is usually between 0.5 and 2. We assume that  $s$  has a uniform probability density function (pdf)  $g(s|a, b)$  where  $a$  and  $b$  define the range of variation. The pdf of the image can be derived by marginalizing  $s$ :

$$g(\mathbf{x}|\boldsymbol{\theta}) = \int_a^b g(\mathbf{x}|\boldsymbol{\theta}, s)g(s|a, b)ds \quad (2)$$

This work was supported by US National Institutes of Health grants R01LM010142 and R01GM095658.

<sup>1</sup>C. Huang is with the Department of Biomedical Engineering, Yale University, New Haven, CT 06520, USA (email: chenxi.huang@yale.edu)

<sup>2</sup>H. D. Tagare is with the Department of Diagnostic Radiology, Electrical Engineering and Biomedical Engineering, Yale University, New Haven, CT 06520, USA (email: hemant.tagare@yale.edu)

where  $g(\mathbf{x}|\boldsymbol{\theta}, s)$  is the pdf of a multivariate normal distribution  $N(s\boldsymbol{\theta}, \sigma^2\mathbf{I})$  based on the inlier model in (1).

### B. Proposed Estimator

We first give the definition of a w-estimate. A w-estimate  $\mathbf{T}$  is the following weighted average of observations  $\mathbf{x}_i, i = 1, 2, \dots, N$ :

$$\mathbf{T} = \frac{\sum \mathbf{x}_i w(\mathbf{x}_i, \mathbf{T})}{\sum w(\mathbf{x}_i, \mathbf{T})} \quad (3)$$

where  $w(\mathbf{x}, \mathbf{T}) \geq 0$  is the weight function that depends on both the observation  $\mathbf{x}$  and the estimate  $\mathbf{T}$ . The estimate is a fixed point of the function defined in (3). The estimate is calculated by the iteration:  $\mathbf{T}^{(j+1)} = \frac{\sum \mathbf{x}_i w(\mathbf{x}_i, \mathbf{T}^{(j)})}{\sum w(\mathbf{x}_i, \mathbf{T}^{(j)})}$  until  $\mathbf{T}^{(j)}$  converges [4]. The converged estimate is denoted as  $\hat{\mathbf{T}}$ .

The properties of the w-estimator depend on the weight function. Our weight function is defined as:

$$w(\mathbf{x}, \boldsymbol{\theta}) = \frac{|\mathbf{x}^\top \boldsymbol{\theta}|}{\|\mathbf{x}\| \|\boldsymbol{\theta}\|} \exp\left[-\beta \left\| \mathbf{x} - \frac{(\mathbf{x}^\top \boldsymbol{\theta})}{\|\boldsymbol{\theta}\|^2} \boldsymbol{\theta} \right\|^2\right] \quad (4)$$

where  $\beta$  is a tuning parameter whose value is discussed later in this paper.

The motivation for the two terms in the right hand side of equation (4) is as follows: The term  $\frac{|\mathbf{x}^\top \boldsymbol{\theta}|}{\|\mathbf{x}\| \|\boldsymbol{\theta}\|}$  is the absolute value of the correlation coefficient of  $\mathbf{x}$  and  $\boldsymbol{\theta}$  and is the term of primary interest. It weighs an image according to its correlation coefficient, and as mentioned before, such weighting is known to be a useful image quality measure for cryo-EM. The second term, i.e.,  $\exp\left[-\beta \left\| \mathbf{x} - \frac{(\mathbf{x}^\top \boldsymbol{\theta})}{\|\boldsymbol{\theta}\|^2} \boldsymbol{\theta} \right\|^2\right]$  depends on the the component of  $\mathbf{x}$  that is orthogonal to  $\boldsymbol{\theta}$ . This term is necessary for the influence function to be bounded.

Asymptotically  $\mathbf{T}$  can be written as a statistical functional of a multivariate distribution  $F$ :

$$\mathbf{T}(F) = \frac{\int \mathbf{x} w(\mathbf{x}, \mathbf{T}(F)) dF(\mathbf{x})}{\int w(\mathbf{x}, \mathbf{T}(F)) dF(\mathbf{x})} \quad (5)$$

$\mathbf{T}(F)$  is Fisher consistent if  $\mathbf{T}(F) = \boldsymbol{\theta}$  where  $dF(\mathbf{x}) = g(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x}$  ( $g(\mathbf{x}|\boldsymbol{\theta})$  is defined in (2)). It is easy to prove that for the choice of  $w$  in (4),  $\mathbf{T}$  is Fisher consistent in the direction, i.e.,  $\mathbf{T} = \alpha \boldsymbol{\theta}$  for some  $\alpha$ . Although the norm of  $\mathbf{T}$  is not consistent, it is easy to show that the difference between  $\|\mathbf{T}\|$  and  $\|\boldsymbol{\theta}\|$  is negligible if  $\|\boldsymbol{\theta}\|/\sigma > 10$  and the pdf  $g(\mathbf{x}|\boldsymbol{\theta})$  is symmetric about  $\boldsymbol{\theta}$  (the proof is omitted due to space limit). These conditions are satisfied for most real cryo-EM images, so we will regard  $\mathbf{T}$  as Fisher consistent for the remainder of the paper.

### C. The Influence Function

The influence function of  $\mathbf{T}$  at  $F$  is given as [4]

$$IF(\mathbf{x}; \mathbf{T}, F) = \lim_{\epsilon \rightarrow 0} \frac{\mathbf{T}(F_\epsilon) - \mathbf{T}(F)}{\epsilon} = \frac{\partial}{\partial \epsilon} [\mathbf{T}(F_\epsilon)]_{\epsilon=0} \quad (6)$$

where  $F_\epsilon = (1 - \epsilon)F + \epsilon \Delta_{\mathbf{x}}$  is a contaminated distribution at the point  $\mathbf{x}$ . IF quantifies asymptotically the influence of the contamination on the estimate.

**Proposition 1** For a Fisher consistent  $\mathbf{T}$  defined in (5) with a differentiable weight function  $w(\mathbf{x}, \boldsymbol{\theta})$ , if  $\|M(\boldsymbol{\theta})\| < 1$  where  $\|\cdot\|$  is the operator norm, then

$$IF(\mathbf{x}; \mathbf{T}, F) = [\mathbf{I} - M(\boldsymbol{\theta})]^{-1} \frac{w(\mathbf{x}, \boldsymbol{\theta})(\mathbf{x} - \boldsymbol{\theta})}{\int w(\mathbf{y}, \boldsymbol{\theta}) dF(\mathbf{y})} \quad (7)$$

where  $M(\boldsymbol{\theta}) = \frac{\int (y - \boldsymbol{\theta}) \frac{\partial}{\partial \boldsymbol{\mu}} [w(\mathbf{y}, \boldsymbol{\mu})]_{\boldsymbol{\mu}=\boldsymbol{\theta}} dF(\mathbf{y})}{\int w(\mathbf{y}, \boldsymbol{\theta}) dF(\mathbf{y})}$ .

The influence function  $IF(\mathbf{x}; \mathbf{T}, F)$  measures the effect of an outlier at point  $\mathbf{x}$  on the estimate  $\mathbf{T}$ . It is desirable for the influence function to be bounded at  $\mathbf{x}$  where  $\mathbf{x}$  is an outlier. As mentioned in the Introduction, outliers in cryo-EM have very low correlation coefficients with the average image and a finite component along the average. A simple way of evaluating influence of such outliers is to consider the value of the influence function as the correlation between  $\mathbf{x}$  and  $\boldsymbol{\theta}$  goes to zeros while the component of  $\mathbf{x}$  along  $\boldsymbol{\theta}$  remains finite. The next proposition states sufficient conditions for the influence function of our w-estimator with weight function to be bounded under these conditions.

**Proposition 2** For  $\mathbf{T}$  and  $w$  defined in (5) and (4), pdf  $g(\mathbf{x}|\boldsymbol{\theta})$  defined in (2) and  $M(\boldsymbol{\theta})$  defined in Proposition 1, an upper bound of  $\|M(\boldsymbol{\theta})\|$  is (for any positive  $u$ )

$$B(\beta, u, \|\boldsymbol{\theta}\|, \sigma) = \frac{Q_1 \Gamma(\frac{p-2}{2})}{Q_2 \gamma(\frac{p-2}{2}, (\bar{\beta} + 0.5)u^2)} \quad (8)$$

where  $\mathbf{x} \in \mathbb{R}^p$ ,  $\bar{\beta} = \beta \sigma^2$ ,

$\gamma$  is the lower incomplete gamma function, and  $Q_1 = \iint \frac{|y_j|}{\|\boldsymbol{\theta}\|} e^{-\bar{\beta} y_j^2} (1 + 2\bar{\beta} y_1 |y_j|) g(y_1 | \|\boldsymbol{\theta}\|) dy_1 dy_j$ ,

$$Q_2 = \iint \frac{|y_1| e^{-(\bar{\beta} \sigma^2 + 0.5) \sum_2^p y_i^2}}{\sqrt{y_1^2 + y_j^2 + u^2}} g(y_1 | \|\boldsymbol{\theta}\|) dy_1 dy_j,$$

where  $\|\boldsymbol{\theta}\| = \|\boldsymbol{\theta}\|/\sigma$ ,  $g(y|\boldsymbol{\theta}) = \int_a^b g(y|\boldsymbol{\theta}, s) g(s)^{1/p} ds$  and  $g(y|\boldsymbol{\theta}, s)$  is the pdf of  $N(s\boldsymbol{\theta}, \sigma^2)$ .

If  $B(\beta, u, \|\boldsymbol{\theta}\|, \sigma) < 1$ , then  $IF(\mathbf{x}; \mathbf{T}, F) = 0$  for outlier  $\mathbf{x}$  with  $\frac{|\mathbf{x}^\top \boldsymbol{\theta}|}{\|\mathbf{x}\| \|\boldsymbol{\theta}\|} \rightarrow 0$  and a finite  $\frac{|\mathbf{x}^\top \boldsymbol{\theta}|}{\|\boldsymbol{\theta}\|}$ .

We can see from Proposition 2 that for a fixed  $u$ , by choosing proper  $\beta$ , the influence function of the proposed estimator goes to zero. This shows that outliers have no effects on the estimation.

## III. RESULTS

We present results of using our estimator on both simulated and real cryo-EM images. We calculate the estimate  $\hat{\mathbf{T}}$  from a set of aligned images by the iteration mentioned in Section II-B and use the pixel-wise median of the images as the initial estimate. The value of parameter  $\beta$  is determined by calculating function  $B(\beta, u, \|\boldsymbol{\theta}\|, \sigma)$  defined in Proposition 2. We find that when  $\|\boldsymbol{\theta}\| > 10\sigma$ ,  $u = 11\|\boldsymbol{\theta}\|/\sigma$  and  $p = 10^4$ , if  $\beta = 10^{-5}/\sigma^2$ , the condition  $B(\beta, u, \|\boldsymbol{\theta}\|, \sigma) < 1$  is satisfied. For both datasets below, we use  $\beta = 10^{-5}/\sigma^2$ .

### A. Simulated Data

Simulated Cryo-EM images are generated by projecting the atomic structure of the 50S ribosomal subunit from the Protein Data Bank (PDB ID:1JJ2) with a simulated water

shell [5]. Images are filtered by a contrast transfer function (CTF) with defocus  $1.3\mu\text{m}$  (CTF models the effect of the microscope). We also apply a scale factor having uniform distribution of  $U(0.5, 1.5)$ . Gaussian white noise is added to generate images with SNR around -10dB.

Inliers are images of the ribosome along a fixed projection direction. Outliers are from two categories: misaligned images and pure noise. Misaligned images are images of projection directions orthogonal to that of the inliers. We use a uniform mixture of these two categories to generate the outliers.

In this experiment, we use 60 inliers and 40 outliers.  $|\theta|/\sigma = 22$  which satisfies the condition of  $|\theta|/\sigma > 10$ . The estimate  $\hat{\mathbf{T}}$  and examples of an inlier and outliers of the two categories are given in Fig. 1. Fig. 2 shows the weights of the images at the last iteration. The weights exhibit the contribution of each image to the final estimate  $\hat{\mathbf{T}}$ . The weights associated with outliers are significantly smaller than the weights associated with inliers. Thus the contribution of outliers to the estimate of the mean is greatly diminished. To compare the quality of our estimate with the classical average, we calculate the mean square error (MSE) and SNR of the estimates from 100 simulations. The results are shown in Table I. The experiment of simulated data demonstrates the robustness and the consistency of our estimator.

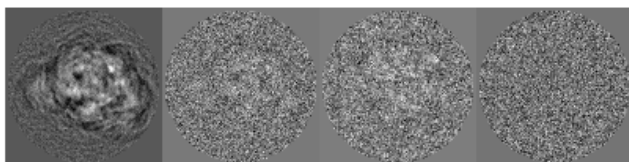


Fig. 1. From left to right (simulated): estimate  $\hat{\mathbf{T}}$ , an inlier, a misaligned image and a pure noise image

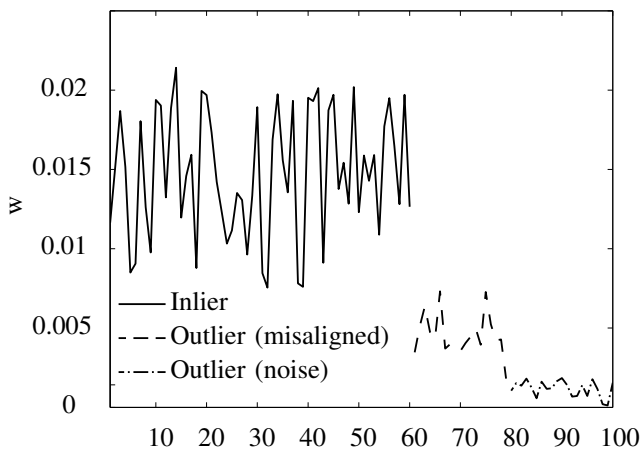


Fig. 2. Weights of simulated images

TABLE I

MEAN SQUARE ERROR (MSE) AND SIGNAL-TO-NOISE RATIO (SNR) OF THE PROPOSED ESTIMATE AND CLASSICAL AVERAGE

	Proposed estimate	Classical average
MSE ( $\times 10^5$ )	1.36	2.25
SNR (dB)	6.7	5.1

## B. Experimental Data

We applied our algorithm to the real 50S ribosomal subunit cryo-EM images that are available from the National Resource for Automated Molecular Microscopy [6]. The images are aligned by the software package SPIDER [7] and we use 60 aligned images along a projection direction. The noise standard deviation is estimated by  $\sigma = \text{median}(|W_{HH}(\mathbf{x})|)/0.6745$ , where  $W_{HH}(\mathbf{x})$  is the wavelet coefficients in the HH (high-high) subband of an image that contains mostly white noise [8]. The estimated noise standard deviation is  $\sigma \approx 0.56$ .  $|\theta|$  is approximated by the norm of the estimate  $|\hat{\mathbf{T}}|$  and we have  $|\theta|/\sigma \approx 40$ . Fig. 3 shows the weights of the images at the last iteration. Two images (4th and 14th) have significant lower weights than the rest of the images. Visualizing these two images (3rd and 4th images in Fig. 4) suggests that they do not appear to contain the signal of the projected structure. Fig.4 also shows the estimate  $\hat{\mathbf{T}}$  and one image that has a high weight.

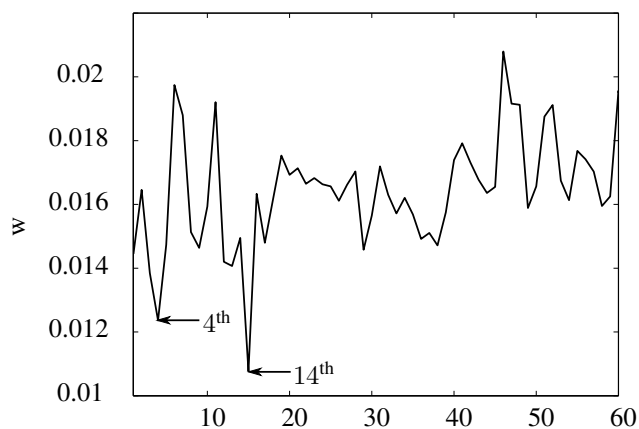


Fig. 3. Weights of cryo-EM images

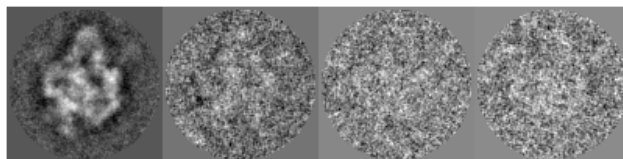


Fig. 4. From left to right (cryo-EM images): estimate  $\hat{\mathbf{T}}$ , an inlier, the 4th image, the 14th image

## IV. CONCLUSIONS

We proposed a method for robustly averaging aligned images for cryo-EM Single Particle Reconstruction. We introduce a novel w-estimator to replace the calculation of average images by a weighted average. The weight function is designed to be asymptotically robust against outliers in cryo-EM and insensitive to signal contrast variation among the inliers. We are able to verify the consistency of our estimator and derive its influence function. We proved that the influence function is bounded. Experiments with simulated data shows good performance on limiting the influence of outlier images on the estimate. Application of our method to real cryo-EM images demonstrates its ability to identify possible outlier images.

## APPENDIX

### A. Proof of Proposition 1

$\mathbf{T}$  defined in (5) at distribution  $F_\epsilon = (1 - \epsilon)F + \epsilon\Delta_{\mathbf{x}}$  satisfies

$$[f w(\mathbf{y}, \mathbf{T}(F_\epsilon)) dF_\epsilon] \mathbf{T}(F_\epsilon) = \int \mathbf{y} w(\mathbf{y}, \mathbf{T}(F_\epsilon)) dF_\epsilon.$$

Taking the derivative with respect to  $\epsilon$  and evaluate at  $\epsilon = 0$ :

$$\begin{aligned} & [f w(\mathbf{y}, \boldsymbol{\theta}) dF + f \boldsymbol{\theta} \frac{\partial}{\partial \boldsymbol{\mu}} [w(\mathbf{y}, \boldsymbol{\mu})]_{\boldsymbol{\theta}} dF] \frac{\partial \mathbf{T}(F_\epsilon)}{\partial \epsilon} \Big|_{\epsilon=0} \\ & + \boldsymbol{\theta} \int w(\mathbf{y}, \boldsymbol{\theta}) d[\Delta_{\mathbf{x}} - F] = \int \mathbf{y} \frac{\partial}{\partial \boldsymbol{\mu}} [w(\mathbf{y}, \boldsymbol{\mu})]_{\boldsymbol{\theta}} dF \frac{\partial \mathbf{T}(F_\epsilon)}{\partial \epsilon} \Big|_{\epsilon=0} \\ & + \int \mathbf{y} w(\mathbf{y}, \boldsymbol{\theta}) d[\Delta_{\mathbf{x}} - F] \end{aligned}$$

$$\begin{aligned} & \Rightarrow [f w(\mathbf{y}, \boldsymbol{\theta}) dF - f(\mathbf{y} - \boldsymbol{\theta}) \frac{\partial}{\partial \boldsymbol{\mu}} [w(\mathbf{y}, \boldsymbol{\mu})]_{\boldsymbol{\theta}} dF] \frac{\partial \mathbf{T}(F_\epsilon)}{\partial \epsilon} \Big|_{\epsilon=0} \\ & = w(\mathbf{x}, \boldsymbol{\theta})(\mathbf{x} - \boldsymbol{\theta}) \end{aligned}$$

$$\text{Let } M(\boldsymbol{\theta}) = \frac{\int (\mathbf{y} - \boldsymbol{\theta}) \frac{\partial}{\partial \boldsymbol{\mu}} [w(\mathbf{y}, \boldsymbol{\mu})]_{\boldsymbol{\theta}} dF}{\int w(\mathbf{y}, \boldsymbol{\theta}) dF}.$$

If  $\|M(\boldsymbol{\theta})\| < 1$  then  $\mathbf{I} - M(\boldsymbol{\theta})$  is invertible. We thus have

$$IF(\mathbf{x}; \mathbf{T}, F) = [\mathbf{I} - M(\boldsymbol{\theta})]^{-1} \frac{w(\mathbf{x}, \boldsymbol{\theta})(\mathbf{x} - \boldsymbol{\theta})}{\int w(\mathbf{y}, \boldsymbol{\theta}) dF}.$$

### B. Proof of Proposition 2

Define a coordinate system where  $\boldsymbol{\theta} = [\theta, \dots, 0]^T$ ,  $\mathbf{y} = [y_1, \dots, y_p]^T$ . First normalize  $\mathbf{y}$  by  $\mathbf{y}/\sigma$ ,  $\boldsymbol{\theta}$  by  $\boldsymbol{\theta}/\sigma$  and  $\beta$  by  $\beta\sigma^2$ . We prove that  $M(\boldsymbol{\theta})$  is a diagonal matrix. Since the denominator of  $M(\boldsymbol{\theta})$  is a scalar, we only need to show that the numerator is a diagonal matrix. Let  $a_{ij}$  denote the  $ij^{\text{th}}$  entry of the numerator of  $M(\boldsymbol{\theta})$ . We have  $\frac{\partial w}{\partial \mu_1} \Big|_{\boldsymbol{\theta}} = 0$  and  $\frac{\partial w}{\partial \mu_j} \Big|_{\boldsymbol{\theta}} = \text{sign}(y_1) \frac{y_j \exp\{-\beta \sum_2^p y_i^2\}}{\|\mathbf{y}\| \|\boldsymbol{\theta}\|} (1 + 2\beta y_1^2)$ ,  $j \neq 1$  ( $w(\mathbf{x}, \mathbf{T})$  is differentiable everywhere but one point where  $\mathbf{x}^T \mathbf{T} = 0$ . In practice we can always replace  $w$  around that point by a smooth function. So we will regard  $w$  as differentiable everywhere for this proof). By the independence of noise, we can write  $g(\mathbf{x}|\boldsymbol{\theta}) = \prod_{1 \leq i \leq p} g(x_i|\theta_i)$ . Then  $a_{ij} = 0$ ,  $\forall i \neq j$  and  $a_{11} = 0$ , i.e.,  $M$  is diagonal. An upper bound of  $a_{jj}$ ,  $j \neq 0$  is

$$\begin{aligned} & \int y_j \frac{\partial w}{\partial \mu_j} \Big|_{\boldsymbol{\theta}} g(\mathbf{y}|\boldsymbol{\theta}) d\mathbf{y} \\ & = \int \text{sign}(y_1) \frac{y_j^2 e^{-\beta \sum_2^p y_i^2}}{\|\mathbf{y}\| \|\boldsymbol{\theta}\|} (1 + 2\beta\sigma^2 y_1^2) g(\mathbf{y}|\boldsymbol{\theta}) d\mathbf{y} \\ & \leq \iint \frac{|y_j|}{|\boldsymbol{\theta}|} e^{-\beta y_j^2} (1 + 2\beta y_1 |y_j|) g(y_1|\theta) dy_1 dy_j \\ & \quad \cdot \int \prod_{i \neq 1, j} g(y_i|\theta_i) dy_i \\ & = C \iint \frac{|y_j|}{|\boldsymbol{\theta}|} e^{-\beta y_j^2} (1 + 2\beta y_1 |y_j|) g(y_1|\theta) dy_1 dy_j \\ & \quad \cdot \frac{\Gamma(\frac{p-2}{2})}{2(\beta + 0.5)^{\frac{p-2}{2}}} \end{aligned}$$

where  $C$  is a constant that only depends on the dimension of  $\mathbf{y}$  and  $\text{sign}(y)$  is the sign function. A lower bound on the denominator of  $M(\boldsymbol{\theta})$  is

$$\int w(\mathbf{y}, \boldsymbol{\theta}) dF = \int \frac{|y_1|}{|\mathbf{y}|} e^{-(\beta+0.5) \sum_2^p y_i^2} g(\mathbf{y}|\boldsymbol{\theta}) d\mathbf{y}$$

$$\begin{aligned} & = C \int_0^\infty \left[ \iint \frac{|y_1| e^{-(\beta+0.5) \sum_2^p y_i^2}}{\sqrt{y_1^2 + y_j^2 + r^2}} g(y_1|\theta) dy_1 dy_j \right] \\ & \quad \cdot r^{p-3} e^{-(\beta+0.5)r^2} dr \\ & \geq C \left[ \iint \frac{|y_1| e^{-(\beta+\frac{1}{2}) \sum_2^p y_i^2}}{\sqrt{y_1^2 + y_j^2 + u^2}} g(y_1|\theta) dy_1 dy_j \right] \\ & \quad \cdot \int_0^u r^{p-3} e^{-(\beta+0.5)r^2} dr \\ & = C \left[ \iint \frac{|y_1| e^{-(\beta+0.5) \sum_2^p y_i^2}}{\sqrt{y_1^2 + y_j^2 + u^2}} g(y_1|\theta) dy_1 dy_j \right] \\ & \quad \cdot \frac{\gamma(\frac{p-2}{2}, (\beta + 0.5)u^2)}{2(\beta + 0.5)^{\frac{p-2}{2}}} \end{aligned}$$

where  $\gamma$  is the lower incomplete gamma function and  $u$  can be any positive value. We thus have an upper bound of the diagonal entries of  $M(\boldsymbol{\theta})$  denoted as  $B(\beta, u, |\boldsymbol{\theta}|, \sigma)$ . If  $\|M(\boldsymbol{\theta})\| \leq B(\beta, u, |\boldsymbol{\theta}|, \sigma) < 1$ , from Proposition 1 the influence function of our estimator is

$$IF(\mathbf{x}; \mathbf{T}, F) = [\mathbf{I} - M(\boldsymbol{\theta})]^{-1} \frac{w(\mathbf{x}, \boldsymbol{\theta})(\mathbf{x} - \boldsymbol{\theta})}{\int w(\mathbf{y}, \boldsymbol{\theta}) dF}.$$

$$\|IF(\mathbf{x}; \mathbf{T}, F)\| \leq \|w(\mathbf{x}, \boldsymbol{\theta})(\mathbf{x} - \boldsymbol{\theta})\| K(\boldsymbol{\theta})$$

where  $K(\boldsymbol{\theta}) \triangleq \|[\mathbf{I} - M(\boldsymbol{\theta})]^{-1} / \int w(\mathbf{y}, \boldsymbol{\theta}) dF\|$  is a constant for a fixed  $\boldsymbol{\theta}$ . We show that  $IF(\mathbf{x}; \mathbf{T}, F) = 0$ , when  $\frac{|\mathbf{x}^T \boldsymbol{\theta}|}{\|\mathbf{x}\| \|\boldsymbol{\theta}\|} \rightarrow 0$  and  $\frac{|\mathbf{x}^T \boldsymbol{\theta}|}{\|\boldsymbol{\theta}\|}$  is finite by showing that  $\|w(\mathbf{x}, \boldsymbol{\theta})(\mathbf{x} - \boldsymbol{\theta})\| \rightarrow 0$ . Let  $x_1$  denote the component of  $\mathbf{x}$  along  $\boldsymbol{\theta}$  and  $x_2$  the component orthogonal to  $\boldsymbol{\theta}$ . Let  $\phi = \frac{|x_1|}{|\mathbf{x}|}$  and  $x_1$  be finite. We thus have

$$\begin{aligned} \lim_{\phi \rightarrow 0} \|w(\mathbf{x}, \boldsymbol{\theta})(\mathbf{x} - \boldsymbol{\theta})\| & = \lim_{\phi \rightarrow 0} \frac{|x_1| e^{-\beta x_2^2}}{\sqrt{x_1^2 + x_2^2}} \sqrt{(x_1 - \theta)^2 + x_2^2} \\ & = \lim_{\phi \rightarrow 0} \phi e^{-\beta(1-\phi^2)x_1^2/\phi^2} \sqrt{\phi^2(x_1 - \theta)^2 + (1-\phi^2)x_2^2/\phi^2} \\ & = \lim_{\phi \rightarrow 0} e^{-\beta x_1^2/\phi^2} |x_1| = 0. \end{aligned}$$

## REFERENCES

- [1] T. R. Shaikh, R. Trujillo, J. S. LeBarron, W. T. Baxter, and J. Frank, "Particle-verification for single-particle, reference-based reconstruction using multivariate data analysis and classification," *Journal of structural biology*, vol. 164, pp. 41–48, Oct. 2008. PMID: 18619547 PMID: PMC2577219.
- [2] W. V. Nicholson and R. M. Glaeser, "Review: automatic particle detection in electron microscopy," *Journal of structural biology*, vol. 133, pp. 90–101, Mar. 2001. PMID: 11472081.
- [3] J. W. Tukey, *Exploratory data analysis*. Reading, Mass.: Addison-Wesley Pub. Co., 1977.
- [4] F. R. Hampel, *Robust statistics: the approach based on influence functions*. New York: Wiley, 1986.
- [5] A. Kucukelbir, F. J. Sigworth, and H. D. Tagare, "A bayesian adaptive basis algorithm for single particle reconstruction," *Journal of Structural Biology*, vol. 179, pp. 56–67, July 2012.
- [6] N. R. Voss, D. Lyumkis, A. Cheng, P.-W. Lau, A. Mulder, G. C. Lander, E. J. Brignole, D. Fellmann, C. Irving, E. L. Jacovetty, A. Leung, J. Pulokas, J. D. Quispe, H. Winkler, C. Yoshioka, B. Carragher, and C. S. Potter, "A toolbox for ab initio 3-d reconstructions in single-particle electron microscopy," *Journal of structural biology*, vol. 169, pp. 389–398, Mar. 2010. PMID: 20018246 PMID: PMC2826578.
- [7] J. Frank, M. Radermacher, P. Penczek, J. Zhu, Y. Li, M. Ladjadj, and A. Leith, "SPIDER and WEB: processing and visualization of images in 3D electron microscopy and related fields," *Journal of structural biology*, vol. 116, pp. 190–199, Feb. 1996. PMID: 8742743.
- [8] D. Donoho, "De-noising by soft-thresholding," *IEEE Transactions on Information Theory*, vol. 41, pp. 613–627, May 1995.