



# Variance estimations in the presence of intermittent interference and their applications to incoherent scatter radar signal processing

Qihou Zhou<sup>1</sup>, Yanlin Li<sup>1</sup>, and Yun Gong<sup>2</sup>

<sup>1</sup>Department of Electrical and Computer Engineering, Miami University, Oxford, Ohio, USA

<sup>2</sup>School of Electronic Information, Faculty of Information Sciences, Wuhan University, Wuhan, China

**Correspondence:** Qihou Zhou (zhouq@miamioh.edu)

Received: 27 January 2024 – Discussion started: 14 February 2024

Revised: 28 April 2024 – Accepted: 26 May 2024 – Published: 17 July 2024

**Abstract.** We discuss robust estimations for the variance of normally distributed random variables in the presence of interference. The robust estimators are based on either ranking or the geometric mean. For the interference models used, estimators based on the geometric mean outperform the rank-based ones in both mitigating the effect of interference and reducing the statistical error when there is no interference. One reason for this is that estimators using the geometric mean do not suffer from the “heavy tail” phenomenon like the rank-based estimators do. The ratio of the standard deviation over the mean of the power random variable is sensitive to interference. It can thus be used as a criterion to combine the sample mean with a robust estimator to form a hybrid estimator. We apply the estimators to the Arecibo incoherent scatter radar signals to determine the total power and Doppler velocities in the ionospheric E-region altitudes. Although all the robust estimators selected deal with light contamination well, the hybrid estimator is most effective in all circumstances. It performs well in suppressing heavy contamination and is as efficient as the sample mean in reducing the statistical error. Accurate incoherent scatter radar measurements, especially at nighttime and at E-region altitudes, can improve studies of ionospheric dynamics and composition.

such applications. In some cases, the variance represents the power of the undesired noise. In other cases, the variance is the desired signal power, such as in our study here on incoherent scatter radar (ISR) signals. Our broad objective is to explore methods that estimate the variance in a normally distributed random variable accurately in the presence of interference. The general problem falls under robust statistics (e.g., Huber and Ronchetti, 2009; Wilcox, 2017). Specifically, we attempt to optimize ISR signal processing using robust estimators.

An ISR, with a large aperture and high transmitting power, measures the electron concentration and other state variables in the ionosphere. Its versatility makes it the most important ground-based instrument for ionospheric studies. Several major ISRs started operation in the 1960s. Readers are referred to Evans (1969) for the principle, capabilities, and comparisons of the early facilities. An ISR typically transmits a binary-phase code to increase the signal-to-noise ratio. The received signal consists of sequences of altitude-dependent in-phase and quadrature voltage samples, which, upon decoding, can be used to obtain a variety of ionosphere parameters such as electron density and electron and ion temperatures (e.g., Zhou et al., 1997; Isham et al., 2000; Hysell et al., 2014). An essential characteristic of the voltage samples is that they are normally distributed, with the variance proportional to the electron density at the corresponding altitude. Because an ISR measures the tiny amount of power scattered off the electrons and ions in space, averaging over 1000 samples is essential to derive ionospheric parameters. In the absence of interference, a simple arithmetic average of the voltage samples squared provides the best estimator for the total power or power spectral density estimates, which

## 1 Introduction

In radar signal processing and in many other applications, the data samples can often be modeled as a constant superimposed on a normally distributed random variable. The variance of the random process is an important parameter in

form the foundation for the derivation of various ionospheric and atmospheric variables. It is well known, however, that the sample mean is susceptible to outliers. In many cases, it is necessary to use other estimators to obtain meaningful averages.

The ISR signal is subject to both active and passive interference. The former can be from other radars and TV stations. The latter can be from scattering off ships, satellites, and other objects. The most significant interference source for ISRs is micro-meteors, although they are the desired signal in the context of meteor studies (e.g., Zhou et al., 1995; Chau et al., 2007; Li et al., 2023). Meteor echoes come in diverse strengths and durations and provide the physical basis for constructing the interference model in our simulations. The incoherent scatter radar signal provides a textbook case for a normally distributed random variable that exists in nature. The high sensitivity of an ISR makes it susceptible to various types of interference. ISR signals thus provide a good test bed to evaluate the performance of various estimators.

In the following section, we discuss the static characteristics of various estimators and compare their performance through theoretical analysis and numerical simulations for different interference scenarios. The aim here is to find an estimator that performs well with and without interference. In Sect. 3, we compare the performance of several estimators for total ISR power and Doppler velocity processing. We show that the hybrid estimator performs best for practically all the interference scenarios, and it is essentially as effective as the sample mean in reducing the statistical error.

## 2 Characteristics and comparison of mean power estimators

### 2.1 Signal and interference models

Let  $X$  be an independent identically distributed normal random variable having  $N = N_1 N_2$  data samples organized as

$$X = \begin{Bmatrix} x_{11} & \cdots & x_{1N_1} \\ \vdots & \ddots & \vdots \\ x_{N_2 1} & \cdots & x_{N_2 N_1} \end{Bmatrix}.$$

For radar and many other digital sampling systems,  $X \sim N(0, \sigma^2)$  can be regarded as voltage samples.  $Y = \left\{ \frac{1}{N_1} \sum_{n_1=1}^{N_1} x_{1n_1}^2, \frac{1}{N_1} \sum_{n_1=1}^{N_1} x_{2n_1}^2, \dots, \frac{1}{N_1} \sum_{n_1=1}^{N_1} x_{N_2 n_1}^2 \right\}$  represents the power random variable with  $N_2$  elements. Each element in  $Y$  is a sample mean of  $N_1$  raw power variables,  $X^2$ . The expectation of  $Y_i$  is  $\sigma_0^2$ , which is the variance of  $X$ . We strive to estimate  $\sigma_0^2$  most accurately, given samples of  $X$ . As there are many types of variances, we will call estimating  $\sigma_0^2$  “power estimation” to be specific and to minimize confusion. In the absence of interference,  $Y_i$  can be shown to have a gamma probability density distribution (pdf):

$$f\left(y; \frac{N_1}{2}, \frac{2\sigma_0^2}{N_1}\right) = \frac{y^{\frac{N_1}{2}-1} e^{-\frac{yN_1}{2\sigma_0^2}}}{\Gamma\left(\frac{N_1}{2}\right)} \left(\frac{N_1}{2\sigma_0^2}\right)^{\frac{N_1}{2}}, \quad (1)$$

where  $\frac{N_1}{2}$  and  $\frac{2\sigma_0^2}{N_1}$  are the shape and scale parameters, respectively, and the support of  $y$  is  $(0, \infty)$  (e.g., Wikipedia, 2024b). The corresponding cumulative distribution function is

$$\begin{aligned} F\left(y; \frac{N_1}{2}, \frac{2\sigma_0^2}{N_1}\right) &= \frac{1}{\Gamma\left(\frac{N_1}{2}\right)} \gamma\left(\frac{N_1}{2}, \frac{N_1}{2\sigma_0^2} y\right) \\ &= \frac{1}{\Gamma\left(\frac{N_1}{2}\right)} \int_0^{\frac{N_1}{2\sigma_0^2} y} t^{\frac{N_1}{2}-1} e^{-t} dt, \end{aligned} \quad (2)$$

where  $\gamma(s, x) = \int_0^x t^{s-1} e^{-t} dt$  is the lower incomplete gamma function. Distribution function  $f(y)$  can also be viewed as a  $N_1$ -degree chi-squared distribution scaled by  $N_1$ . The variance of  $Y_i$  is  $\frac{2\sigma_0^4}{N_1}$ . The distribution functions at  $N_1 = 1, 2,$  and  $8,$  which we will study in more detail, are  $f(y; \frac{1}{2}, 2\sigma_0^2) = \frac{e^{-\frac{y}{2\sigma_0^2}}}{\sqrt{2\pi y \sigma_0^2}}, f(y; 1, \sigma_0^2) = \frac{e^{-\frac{y}{\sigma_0^2}}}{\sigma_0^2},$  and  $f(y; 4, \frac{\sigma_0^2}{4}) = \frac{2^7 y^3 e^{-\frac{4y}{\sigma_0^2}}}{3\sigma_0^8},$  respectively. At large  $N_1,$  the pdf is approximately normal,  $f(y; \frac{N_1}{2}, \frac{2\sigma_0^2}{N_1}) \sim N\left(\sigma_0^2, \frac{2\sigma_0^4}{N_1}\right).$  Of particular interest is the case of  $N_1 = 2,$  which corresponds to the in-phase and quadrature samples in a radar system.

The interference is also modeled as a gamma distribution with a shape parameter of  $k = 4$  and scale parameter  $(a_\eta \sigma_0)^2/k,$  which has a mean of  $a_\eta^2 \sigma_0^2.$  Since we are mainly concerned with the signal shape parameter being  $1/2$  and  $1,$  a larger shape parameter in the interference model makes it easier to differentiate between interference and signal, as the interference is more concentrated around a higher mean value. The interference is equally likely to occur at each data point, with a probability of  $p_\eta = 0.01,$  and is always additive to the signal. The total interference power relative to the signal power is thus  $p_\eta a_\eta^2.$  We will mainly consider three cases of interference with  $a_\eta = 2, 6,$  and  $18$  to represent low, moderate, and strong interference, respectively.

### 2.2 Estimators and their characteristics in the absence of interference

The most common estimators are the sample mean, geometric mean, and median. The sample mean of  $Y$  is the arithmetic average of  $N_2$  samples, i.e.,  $A_N \equiv \frac{1}{N_2} \sum_{i=1}^{N_2} Y_i,$  where  $Y_i$  is the sample mean of  $X^2$  averaged over  $N_1$  samples. With a known shape parameter, the sample mean is the uniformly minimum-variance unbiased estimator (UMVUE) and maximum likelihood estimator (e.g., Siegrist, 2022; Wikipedia,

2024a). The geometric mean,  $G_N \equiv \left(\prod_{i=1}^{N_2} Y_i\right)^{1/N_2}$ , and median,  $D_N \equiv \text{med}(Y_1, \dots, Y_{N_2})$ , are more resistant to outliers but not effective in reducing the statistical fluctuations. Although the three basic estimators are largely at the opposite ends of efficiency vs. robustness, they can serve as building blocks for other estimators. In the following, we discuss the three basic estimators and compare them with a weighted mean, a hybrid estimator, and two trimmed estimators.

The effectiveness of a power estimator,  $Z$ , in reducing the statistical fluctuation is measured by the normalized variance

$$R^2(Z) \equiv \frac{N\sigma_z^2}{2\mu_z^2}, \tag{3}$$

where  $\sigma_z^2$  and  $\mu_z$  are the variance and mean of the power estimator while the absolute error is of importance in some cases as well. For the sample mean estimator,  $A_N$ , its distribution is expressed by Eq. (1) with  $N_1$  replaced by  $N$ .  $E(A_N)$  is  $\sigma_0^2$  and the variance is  $2\sigma_0^4/N$ . The theoretical expectation of  $R^2(A_N)$  is thus 1 for the sample mean, which is the lowest that one can obtain. The inverse of  $R^2(Z)$  is the efficiency of the estimator. It is of interest to note that since  $N$  averages can be expressed as the weighted means of  $N_1$  and  $N - N_1$  samples, it follows that the convolution of two gamma distributions remains a gamma distribution. This convolution invariance property is also true of most commonly used distributions, including binomial, Poisson, normal, and chi-squared distributions. In general, if the distribution function of the sum or mean remains the same type for different numbers of samples, it is convolution invariant.

The median and its variance do not appear to have a closed form for  $N_1$  and  $N_2$  in general, although there are closed forms for specific  $N_1$  and large  $N$ . Here we derive the theoretical results for  $N_1 = 1, 2, 8$ , and large  $N$ . For large  $N_2$  and an ascending ranking order  $K$  relatively close to  $N_2/2$ , Zhou et al. (1999) show that ranking has an asymptotic normal distribution, with a variance of  $\sigma_{N_2K}^2 = \frac{K(N_2-K)}{N_2^3 f^2(\mu_r)}$ , where  $\mu_r$  is the ranking value (e.g.,  $K = N_2/2$  for median) and  $f(\mu_r)$  is the pdf for the rank random variable, i.e., Eq. (1) for our study here. For the median estimator, the normalized variance is

$$R^2(D_N) = \frac{N_1}{8f^2(\mu_r; N_1/2, 2/N_1)\mu_r^2}. \tag{4}$$

The median can be solved from  $F(\mu_r) = 1/2$ . For  $N_1 = 1$ , the median is  $2\text{ierf}^2(\frac{1}{2})\sigma_0^2 = 0.4549\sigma_0^2$ , where  $\text{ierf}$  is the inverse of the error function  $\frac{2}{\sqrt{\pi}}\int_0^x e^{-t^2} dt$ . For  $N_1 = 2$ , the median is  $\mu_r = \sigma_0^2 \ln 2 = 0.6931\sigma_0^2$ . The median for  $N_1 = 8$  is  $0.9180\sigma_0^2$ , which can be solved from  $\gamma(4, 4\mu_r) = 3$ . For large  $N_1$ , the pdf tends toward normal and the median tends toward  $\sigma_0^2$ . The  $R^2(D_N)$  values for  $N_1 = 1, 2, 8$ , and  $100$  and  $N = 10000$  are 2.7206, 2.0814, 1.6848, and 1.5760, respectively. In the limiting case of  $N_1$  and  $N_2$  tending toward infinity,  $R^2(D_N) = \pi/2$ , indicating that it takes  $\pi/2$  times the

number of samples for the median operator to achieve the same error as the sample mean. Zhou et al. (1999) also show that taking the 79.7% largest value gives the smallest  $R^2$  at 1.5432. (in Zhou et al., 1999,  $\pi/2$  in Eqs. 24 and 26 should have been  $2/\pi$ .)

In Table 1, we list the  $R^2$  values and the absolute errors for eight estimators in the null-interference case. The second column is the mean of each estimator without scaling for  $\sigma_0 = 1$  (the mean is proportional to  $\sigma_0^2$ ). To compare the different estimators on the same scale, the mean is divided by the respective estimator so that all the estimators in all the cases have a mean of 1 for all subsequent computations in the other columns in Tables 1 and 2. The values not in parentheses listed in the two tables are at least 100 000 Monte Carlo simulations with  $N = 10000$  for all estimators except  $H_N$ . The values in parentheses in Table 1 are theoretical predictions that we can derive.

The mean and variance of the geometric mean ( $G_N$ ) can be obtained by first finding the expectation and variance of one element,  $Y_i^{1/N_2}$ , in the product. The expectation of  $Y_i^{1/N_2}$  is

$$E\left(y^{\frac{1}{N_2}}\right) = \int_0^\infty y^{\frac{1}{N_2}} f(y) dy = \frac{\Gamma\left(\frac{N_1}{2} + \frac{1}{N_2}\right)}{\Gamma\left(\frac{N_1}{2}\right)} \left(\frac{2\sigma_0^2}{N_1}\right)^{\frac{1}{N_2}}. \tag{5}$$

The second moment of  $Y_i^{1/N_2}$  is

$$E\left(y^{\frac{2}{N_2}}\right) = \int_0^\infty y^{\frac{2}{N_2}} f(y) dy = \frac{\Gamma\left(\frac{N_1}{2} + \frac{2}{N_2}\right)}{\Gamma\left(\frac{N_1}{2}\right)} \left(\frac{4\sigma_0^4}{N_1^2}\right)^{\frac{1}{N_2}}. \tag{6}$$

Assuming that  $Y_i$ 's are independent, the expectation, second moment, and variance of the geometric mean are, respectively,

$$E(G_N) = \left(E\left(y^{\frac{1}{N_2}}\right)\right)^{N_2} = \frac{\Gamma^{N_2}\left(\frac{N_1}{2} + \frac{1}{N_2}\right)}{\Gamma^{N_2}\left(\frac{N_1}{2}\right)} \frac{2\sigma_0^2}{N_1}, \tag{7}$$

$$E(G_N^2) = \left(E\left(y^{\frac{2}{N_2}}\right)\right)^{N_2} = \frac{\Gamma^{N_2}\left(\frac{N_1}{2} + \frac{2}{N_2}\right)}{\Gamma^{N_2}\left(\frac{N_1}{2}\right)} \frac{4\sigma_0^4}{N_1^2}, \tag{8}$$

$$\text{Var}(G_N) = \frac{4\sigma_0^4}{N_1^2} \left( \frac{\Gamma^{N_2}\left(\frac{N_1}{2} + \frac{2}{N_2}\right)}{\Gamma^{N_2}\left(\frac{N_1}{2}\right)} - \frac{\Gamma^{2N_2}\left(\frac{N_1}{2} + \frac{1}{N_2}\right)}{\Gamma^{2N_2}\left(\frac{N_1}{2}\right)} \right). \tag{9}$$

The normalized variance for the geometric mean,  $R^2(G_N)$ , is thus

$$R^2(G_N) = \frac{N_1 N_2}{2} \left[ \left( \frac{\Gamma\left(\frac{N_1}{2}\right)\Gamma\left(\frac{N_1}{2} + \frac{2}{N_2}\right)}{\Gamma^2\left(\frac{N_1}{2} + \frac{1}{N_2}\right)} \right)^{N_2} - 1 \right]. \tag{10}$$

This equation is precise for all  $N_1$  and  $N_2$ .  $E(G_N)$  and  $R^2(G_N)$  values for  $N = 10000$  and  $N_1 = 1, 2, 8$ , and  $100$  are listed in Table 1. We are not aware of a precise distribution function for  $G_N$  in general. For the asymptotic case of

**Table 1.** Monte Carlo simulations and theoretical values (in parenthesis) of the mean,  $R^2$ , and absolute error for eight estimators when there is no interference.

Method	$N_1$	Mean (theory)	$R^2$	Error
$A_N$ $N = 10000$	1	1.0000 (1)	1.0020 (1)	0.0113
	2	1.0000 (1)	0.9994 (1)	0.0112
	8	1.0000 (1)	1.0077 (1)	0.0113
	100	1.0000 (1)	0.9945 (1)	0.0113
$D_N$ $N = 10000$	1	0.4549; (0.4549)	2.7149; (2.7206)	0.0186
	2	0.6930 (ln 2);	2.0927; (2.0814)	0.0162
	8	0.9176 (0.9180)	1.6980; (1.6848)	0.0147
	100	0.9917 (1)	1.5614; (1.5760)	0.0150
$G_N$ $N = 10000$	1	0.2808; (0.2808)	2.4841; (2.4672)	0.0178
	2	0.5615; (0.5616)	1.6487; (1.6447)	0.0144
	8	0.8780; (0.8780)	1.1377; (1.1352)	0.0120
	100	0.9901; (0.9901)	1.0028; (1.0100)	0.0114
$T_{05}$ $N = 10000$	1	0.7589; (0.7590)	1.1480; (1.1423)	0.0121
	2	0.8424; (0.8430)	1.0901; (1.0898)	0.0117
	8	0.9317; (0.9320)	1.0434; (1.0431)	0.0116
	100	0.9839; (0.9835)	1.0198; (1.0178)	0.0114
$T_{MAD8}$ $N = 10000$	1	0.8742	1.6973	0.0147
	2	0.8425	1.0769	0.0117
	8	1.0000	1.0075	0.0113
	100	1.0000	0.9984	0.0113
$T_{GEO4}$ $N = 10000$	1	0.9979	1.0185	0.0114
	2	0.9884	1.0763	0.0117
	8	0.9987	1.0210	0.0114
	100	1.0000	0.9984	0.0113
$W_N$ $N = 10000$	1	0.9576	1.0419	0.0115
	2	0.9563	1.0431	0.0115
	8	0.9888	1.0167	0.0114
	100	1.0000	0.9995	0.0113
$H_N$ $N = 1000$	1	0.9576	1.0102	0.0360
	2	0.9563	1.0178	0.0356
	8	0.9888	1.0001	0.0357
	100	1.0000	1.0052	0.0358

large  $N_2$ , Zhou et al. (1999) show that the geometric mean tends toward the normal distribution, with the variance as

$$\text{Var}(G_N)|_{N_2 \rightarrow \infty} = \frac{E^2(G_N)\sigma_{\ln}^2}{N_2}, \tag{11}$$

where  $\sigma_{\ln}^2$  is the variance of  $\ln(y)$ .  $\sigma_{\ln}^2$  is known to equal the trigamma function ( $\psi_1(\frac{N_1}{2})$ ; e.g., Wikipedia, 2024a). Thus,

$$R^2(G_N)|_{N_2 \rightarrow \infty} = \frac{N_1\sigma_{\ln}^2}{2} = \frac{N_1}{2}\psi_1\left(\frac{N_1}{2}\right). \tag{12}$$

For the trigamma function,  $\psi_1(\frac{1}{2}) = \frac{\pi^2}{2}$ ,  $\psi_1(1) = \frac{\pi^2}{6}$ , and other  $\psi_1(\frac{N_1}{2})$  values can be found from the re-

currence relation  $\psi_1(z + 1) = \psi_1(z) - 1/z^2$ . The asymptotic  $R^2(G_N)$  for  $N_1 = [1, 2, 8, 100]$  and  $N = 10000$  is  $[2.4674, 1.6449, 1.3529, 1.010]$ , respectively. They are accurate to the third decimal place compared to the exact values obtained from Eq. (10) for  $N_2 = 10000$ . For large  $N_1$  and  $N_2$ ,  $R^2(G_N) \sim 1 + \frac{1}{N_1}$ , which gives the number of initial averages,  $N_1$ , needed to achieve a certain level of efficiency for the geometric mean. The expectation of  $G_N$  for large  $N_2$  is found to be

$$\begin{aligned} E(G_N)|_{N_2 \rightarrow \infty} &\sim \sigma_0^2 \left(1 + \frac{2}{N_1 N_2}\right) e^{-\frac{1}{N_1} - \frac{1}{3(N_1+2/N_2)N_1}} \\ &\sim \sigma_0^2 e^{-\frac{1}{N_1} - \frac{1}{3N_1^2}} \end{aligned} \tag{13}$$

**Table 2.** Mean and  $R^2$  values for low, moderate, and strong interference. The interference occurrence rate is  $p_\eta = 0.01$  for all three interference scenarios.

Method	$N_1$	$a_\eta = 2$		$a_\eta = 6$		$a_\eta = 18$		Avg $R^2$
		Mean	$R^2$	Mean	$R^2$	Mean	$R^2$	
$A_N$ $N = 10000$	1	1.0400	1.0166	1.3599	4.8437	4.2393	36.507	14.122
	2	1.0400	1.0091	1.3600	4.8794	4.2404	36.127	14.005
	8	1.0400	1.0136	1.3601	4.8697	4.2391	36.015	13.966
	100	1.0400	1.0193	1.3600	4.8984	4.2396	36.411	14.110
$D_N$ $N = 10000$	1	1.0237	2.7454	1.0239	2.7388	1.0238	2.7317	2.7486
	2	1.0290	2.1082	1.0294	2.1277	1.0295	2.1273	2.1211
	8	1.0394	1.7302	1.0555	1.9460	1.0554	1.9503	1.8755
	100	1.040	1.6002	1.2779	8.9737	3.4042	111.64	40.738
$G_N$ $N = 10000$	1	1.0278	2.4724	1.0488	2.5614	1.0717	1.6730	2.2356
	2	1.0316	1.6650	1.0701	1.8519	1.1167	2.2247	1.9139
	8	1.0375	1.1544	1.1467	2.0017	1.3376	5.1570	2.7710
	100	1.0399	1.0310	1.3147	4.3308	2.9238	42.091	15.818
$T_{95}$ $N = 10000$	1	1.0358	1.1711	1.0430	1.2202	1.0430	1.2239	1.2051
	2	1.0380	1.1079	1.0561	1.2533	1.0562	1.2498	1.2037
	8	1.0394	1.0594	1.1450	3.5711	1.7359	94.667	33.099
	100	1.0400	1.0394	1.3243	4.8256	3.7994	41.530	15.798
$T_{MAD8}$ $N = 10000$	1	1.0284	1.7843	1.0080	1.6774	1.0080	1.6714	1.7110
	2	1.0398	1.0910	1.0056	1.1300	1.0020	1.0801	1.1004
	8	1.0399	1.0143	1.1085	3.1179	1.0003	1.0893	1.7405
	100	1.0400	1.0193	1.3601	4.9145	4.2330	39.909	15.281
$T_{GEO4}$ $N = 10000$	1	1.0380	1.0343	1.0087	1.1344	0.9993	1.0159	1.0615
	2	1.0280	1.0885	0.9996	1.1170	0.9981	1.0328	1.0794
	8	1.0400	1.0287	1.0758	2.4126	1.0032	1.1384	1.5266
	100	1.0400	1.0310	1.3600	4.9090	3.9775	54.270	20.070
$W_N$ $N = 10000$	1	1.0390	1.0625	1.0074	1.1054	1.0001	1.0429	1.0703
	2	1.0400	1.0558	1.0115	1.1304	1.0098	1.0415	1.0759
	8	1.0400	1.0246	1.1078	3.1480	1.0001	1.0996	1.7574
	100	1.0391	1.0223	1.3501	4.9301	4.0907	41.158	15.703
$H_N$ $N = 1000$	1	1.0392	1.0236	1.0162	1.1090	1.0092	1.0462	1.0596
	2	1.0377	1.0447	1.0124	1.1290	1.0112	1.0447	1.0728
	8	1.0407	1.0199	1.1247	3.3169	1.0101	1.1043	1.8206
	100	1.0409	1.0231	1.3648	4.9272	4.1395	41.780	15.910

using the approximation  $\ln(\Gamma(z)) \sim z \ln(z) - z - \frac{1}{2} \ln(z) + \frac{1}{12z} + \frac{1}{2} \ln(2\pi)$  (Wikipedia, 2024b). The variance of  $G_N$  at large  $N_2$  is

$$\begin{aligned} \text{Var}(G_N) \Big|_{N_2 \rightarrow \infty} &= \frac{E^2(G_N)\sigma_{\ln}^2}{N_2} \\ &= \frac{\psi_1\left(\frac{N_1}{2}\right)}{N_2} \sigma_0^4 e^{-\frac{2}{N_1} - \frac{2}{3N_1^2}}. \end{aligned} \tag{14}$$

In Table 1, we list the theoretical values of the geometric mean and  $R^2$  and their comparisons with the simulated values. We see that the theoretical values agree with simulations very well for all three basic estimators in the various scenarios.

As the median and other ranks are not efficient in reducing the statistical fluctuation, one can average the data within a certain percentile range, which is known as the trimmed or truncated mean. Since interference is additive, we will only be concerned with one-sided trimming below a fraction of  $\beta$ . Let  $b$  be the integer value of  $\beta N_2$ . The trimmed mean at  $\beta$  is  $T_\beta \equiv \frac{1}{b} \sum_{j=1}^b \text{sort}(Y_i)_j$ , where  $\text{sort}(Y_i)$  is  $Y_i$  sorted into ascending order. Let  $F(y_\beta) = \beta$ ,  $\mu_\beta = \frac{1}{\beta} \int_0^{y_\beta} y f(y) dy$  and  $\sigma_\beta^2 = \frac{1}{\beta} \int_0^{y_\beta} y^2 f(y) dy - \mu_\beta^2$ . Stigler (1973) shows that the asymptotic mean and variance of  $T_\beta$  for large  $N_2$  is

$$E(T_\beta) = \mu_\beta \text{ and } \sigma_T^2 = \frac{\left[ \frac{\sigma_\beta^2}{\beta} + \frac{1-\beta}{\beta} (y_\beta - \mu_\beta)^2 \right]}{N_2}, \text{ respectively. The}$$

normalized variance for the trimmed mean is thus

$$R^2(T_\beta) = N_1 \frac{\sigma_\beta^2 + (1 - \beta)(y_\beta - \mu_\beta)^2}{2\beta\mu_\beta^2}. \quad (15)$$

In the following examples,  $N = 10000$ ,  $\beta = 0.95$ , and  $\sigma_0 = 1$ . If  $N_1 = 2$ , then  $y_\beta = 3.843$ ,  $\mu_\beta = 0.7590$ ,  $\sigma_\beta^2 = 0.7747$ , and  $R^2(T_{95}) = 1.1423$ . If  $N_1 = 8$ , we have  $y_\beta = 2.995$ ,  $\mu_\beta = 0.8422$ ,  $\sigma_\beta^2 = 0.5027$ , and  $R^2(T_{95}) = 1.0898$ . When  $N_1$  is 8, then  $y_\beta = 1.9384$ ,  $\mu_\beta = 0.9320$ ,  $\sigma_\beta^2 = 0.1645$ , and  $R^2(T_{95}) = 1.0431$ . For  $N_1 = 100$ , we have  $y_\beta = 1.2435$ ,  $\mu_\beta = 0.9835$ ,  $\sigma_\beta^2 = 0.0153$ , and  $R^2(T_{95}) = 1.0178$ . As seen in Table 1, the  $R_2$  values agree with the simulation very well. It is of interest to note that  $R_2$  is not  $1/0.95 = 1.05$  as intuition might suggest. It varies from 1.142 at  $N_1 = 1$  to 1.018 at  $N_1 = 100$ . When  $N_1 = 1$ , the tail is long and has more variations, leading to a large  $R^2$  value – a tail-wagging-the-dog situation. More averaging makes the tail more stable and  $R^2$  smaller. This phenomenon, and the effects of a “fat tail” or heavy tail, are extensively discussed by Resnick (2007) and Taleb (2022).

To estimate a parameter robustly, we can attempt to identify outliers and exclude them from the average. Most of the outlier classifying methods involve estimating a nominal deviation and using it in a threshold to detect outliers. The median absolute deviation (MAD), defined as  $\text{MAD} = \text{med}(|Y_i - \text{med}(Y)|)$ , is most frequently used to detect outliers (Huber and Ronchetti, 2009). Since only a small fraction of the ISR data is contaminated most of the time, we will classify a data point having 8 MADs above the median as an outlier. The sample mean of all non-outlier points is referred to as the  $T_{\text{MAD}8}$  estimator. When there is no interference,  $R^2(T_{\text{MAD}8})$  is 1.6973, 1.0769, 1.0075, and 0.9984 for  $N_1 = 1, 2, 8$ , and 100, respectively. There is a significant improvement in  $R^2$  from  $N_1 = 1$  to  $N_1 = 8$  because averaging reduces the number of spurious outliers significantly in the trimmed mean as discussed above. At  $N_1 = 1$ , the proportion of flagged outliers is about 0.2 %, while at  $N_1 = 8$  the effective rate of flagged outliers is 0.0012 %. We note that Rousseeuw and Croux (1993) present two robust estimators that are more efficient than MAD, although more computationally intensive. With a normalized variance larger than 1.2, their estimators are better suited for heavy contamination.

As the geometric mean is resistant to outliers as well, it may also conceivably be used to classify outliers. We define the geometric deviation as  $\sigma_G \equiv G_N e^{\sigma_{\log}(y)} - G_N$ , where  $\sigma_{\log}(y) = \text{std}(\ln(y))$ . The dimensionless  $e^{\sigma_{\log}(y)}$  is known as the geometric standard deviation.  $\sigma_G$  is zero if all samples in  $Y$  are a constant and increase in proportion with  $Y$ , although  $\sigma_G^2$  does not have the usual properties of the variance as commonly defined. We average all the data points 4 geometric deviations below the geometric mean and refer to the estimator as  $T_{\text{GEO}4}$ .  $T_{\text{GEO}4}$  and  $T_{\text{MAD}8}$  are chosen to have almost the same normalized variance at  $N_1 = 2$ , as they flag

out the same number of outliers in the absence of interference. When  $N_1 = 1$ ,  $T_{\text{GEO}4}$  has a far better  $R^2$  value in the null-interference case.

Weighted means can also be used to mitigate the effect of outliers and interference. In this method, values far away from the expected mean are weighted less than those points around the mean. The weighting function we choose is  $w_i = \frac{(y_i - m_{G4})^2}{40\sigma_{G4}^2}$ , where  $m_{G4}$  and  $\sigma_{G4}$  are the sample mean and standard deviation of the  $T_{\text{GEO}4}$  estimator discussed above. The mean values of  $W_N$  for various  $N_1$  are listed in Table 1. In the null-interference case,  $R^2(W_N)$  is no larger than 1.046, or the efficiency is no less than 95.6 %. If the constant 40 is changed to 60, the worst  $R^2$  becomes 1.031, but the weighted mean is less effective in mitigating the effect of outliers. The mean and standard deviation of  $T_{\text{GEO}4}$  are chosen because of their general accuracy and computing efficiency.

Knowing whether interference exists can help mitigate its effect. For a gamma distribution, we cannot associate the existence of outliers with interference with certainty, as there are outliers even when there is no interference. Since the expectation of  $R^2$  for the sample mean is known in the null-interference case, a deviation from the expectation indicates that the underlying process may contain interference. As the sample mean performs best when there is no interference, an expedient strategy to reduce the variance is to combine the sample mean when no interference is detected with another estimator that is effective in mitigating the interference. We have used  $N = 10000$  for the asymptotic case for all the estimators discussed above. In combining different estimators, a smaller  $N$  value is preferred so that the combined estimator will not be dominated by the interference-mitigating estimator in the presence of interference. We can also define a mixed  $R^2$  that uses the mean of the  $T_{\text{GEO}4}$  estimator and the normally defined variance. Such a mixed  $R^2$  is more sensitive to outliers, but its variance is larger. Simulations show it does not cause a material difference from  $R^2(A_N)$  using the sample mean and standard deviation. Because of its simplicity, we choose  $R^2(A_N)$  as the criterion to determine if the data samples follow the desired process. The decision rule for this hybrid estimator,  $H_N$ , is that if  $R(A_N)$  is less than 2 standard deviations above the mean, it uses the sample mean, otherwise the weighted mean is used. The performance of such a combined or hybrid estimator compares well to the other estimators. In Table 1,  $N$  is 1000 for the hybrid estimator,  $H_N$ .

As seen in Table 1, all the order-based estimators ( $D_N$ ,  $T_{95}$ , and  $T_{\text{MAD}8}$ ) perform better as  $N_1$  increases. The tail-wagging-the-dog phenomenon discussed for  $T_{95}$  above is also applicable to  $D_N$  and  $T_{\text{MAD}8}$ , as they also truncate the largest values. Although  $T_{\text{GEO}4}$  is also a trimmed mean, the tail does not control  $R^2$  in the same manner as in the order-based estimators because the length of the tail depends on the largest values. Large sample values increase the geometric deviation, which diminishes the chance of a large sam-

ple value being counted as an outlier. Compared to  $T_{\text{MAD8}}$ ,  $T_{\text{GEO4}}$  flags out fewer outliers at  $N_1 = 1$  but more outliers at  $N_1 = 8$ . At very large  $N_1$  (e.g., 100), the pdf of  $Y_i$  is approximately normal, and all the estimators perform equally well at the theoretical best. It is of interest to note that  $R^2(W_N)$  for the weighted mean is not a strong function of  $N_1$ . The hybrid estimator  $R^2$  is always less than 1.02, making the efficiency better than 98 % for all  $N_1$ 's when there is no interference.

### 2.3 Comparison of estimators in the presence of interference

In Table 2, we list the mean and  $R^2$  values with three levels of noise for the eight estimators discussed above. The total noise power is the mean of  $A_N$  subtracted from 1, which is set as the signal power. In the low-noise case,  $a_\eta = 2$ , the total noise power is 4 % of the signal power. We see that the expectation of the sample mean is 1.04 irrespective of  $N_1$ , as the power is additive. In this case of low interference power, the performance of all the estimators does not differ from the null-interference case significantly. For moderate- and high-noise cases, all the estimators perform very poorly at  $N_1 = 100$ , as practically all the  $Y_i$ 's are contaminated.  $T_{\text{MAD8}}$  performs the best at  $N_1 = 8$  and 100 for  $a_\eta = 18$ . In general, rank-based estimators do better than geometric-mean-based estimators when a large portion of data is contaminated. Large  $N_1$  is akin to having a higher percentage of interference and therefore should be avoided. The strong interference case is easier to deal with than the moderate case is, as it has a very distinct distribution from the signal distribution. The most challenging case is the moderate interference case,  $a_\eta = 6$ . All the estimators perform worse than in the other two interference scenarios. For the moderate case of interference, the weighted mean performs the best at  $N_1 = 1$ , while  $T_{\text{GEO4}}$  does the best at  $N_1 = 2$ .

The last three robust estimators, all of which are based on the geometric mean, have about the same performance. They perform better than the rank-based estimators at  $N_1 = 1$  and 2. The averages of the  $R^2$  values for the three noise levels are listed in the last column in Table 2. On balance, the hybrid estimator performs best for the two cases of small  $N_1$ . It should be noted that simulations for the hybrid estimator are based on  $N = 1000$  in Table 2 but on  $N = 10000$  for other estimators. It is almost certain that the hybrid estimator performs the same as  $W_N$  does at modest and strong interference. At low interference levels,  $H_N$  outperforms  $W_N$  because of the inclusion of the sample mean. Thus, the hybrid estimator combining  $W_N$  and  $A_N$  would always perform better than  $W_N$ . The reason that  $R^2(H_N)$  is not always smaller than  $R^2(W_N)$  in some cases in Table 2 is because the statistics at  $N = 1000$  are slightly inferior to those at  $N = 10000$ . Similarly, an estimator combining  $T_{\text{GEO4}}$  with  $A_N$  will outperform  $T_{\text{GEO4}}$  for the same  $N$ . Although the performance of the estimators will change if the underlying assumptions are changed,  $H_N$ ,  $T_{\text{GEO4}}$ , and  $W_N$  are the preferred estimators

because of their interference-mitigating ability, efficiency in reducing statistical fluctuation, and computational efficiency. When  $p_\eta$  is less than 0.005,  $W_N$  (by extension, the combination of  $W_N$  and  $A_N$ ) outperforms  $T_{\text{GEO4}}$  for all interference levels. In cases of prevalent contamination (e.g.,  $p_\eta > 10\%$ ) one can combine order-based estimators (such as the median or trimmed mean) with the sample mean.

## 3 Application to incoherent scatter radar signal processing

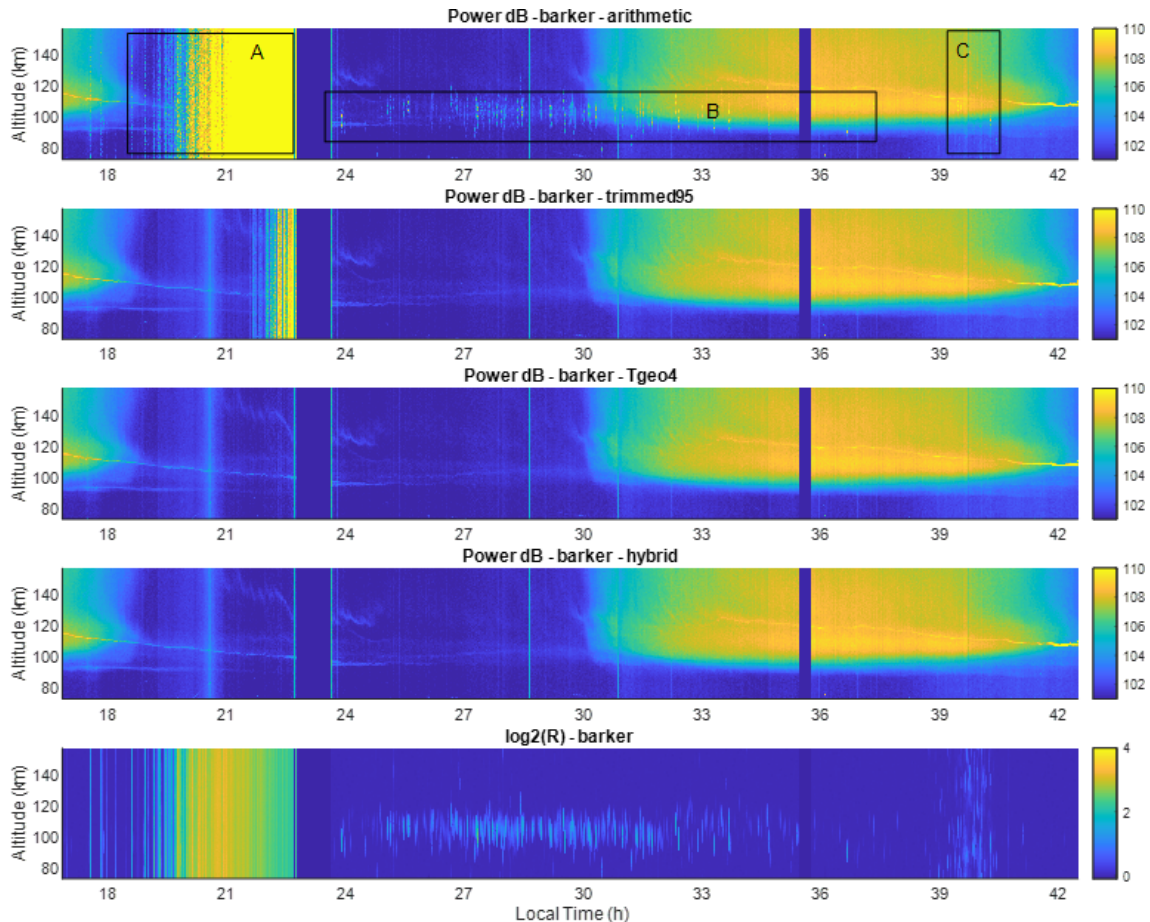
In this section, we apply four estimators to incoherent scatter total power and Doppler velocity processing and compare their performance. The example incoherent scatter radar data were taken at the Arecibo Observatory, Puerto Rico, on 11–12 September 2014. The total power is used to derive the electron density. The Doppler velocity is the same as the neutral wind velocity below about 115 km, but it also depends on the electric field and ion-neutral collision frequency above this altitude. Readers are referred to Zhou et al. (1997) and Isham et al. (2000) for further description of the Arecibo ISR, especially concerning E-region signal processing.

### 3.1 Total power processing

The most common way to obtain the total power and hence electron density in the ionosphere using an ISR is to transmit a 13-baud Barker code with a total pulse length duration less than 52  $\mu\text{s}$ . Barker code is chosen because of its minimized sidelobe. The lack of longer Barker codes is not a severe limitation due to the finite correlation time of the ionosphere. The 13-baud Barker data we use here have a baud length of 2  $\mu\text{s}$ , making the range resolution 300 m. In-phase and quadrature voltage samples from each pulse are stored for post-processing. An inter-pulse period of 10 ms was used so that range aliasing is negligible. As the antenna was pointing vertically, range and altitude are interchangeable here. Although the sampling range in the data was from 60 to 766 km, we mostly focus on the altitude range from 90 to 150 km, where interference is most severe. The raw voltage samples were decoded using a matched filter.

Figure 1 shows the averaged power returns as a function of time and altitude using the sample, trimmed,  $T_{\text{GEO4}}$ , and hybrid means. Because the radar samples are in in-phase and quadrature pairs and larger  $N_1$  contaminates more data samples,  $N_1$  is chosen to be 2. The last panel shows the normalized standard deviation  $R(A_N)$  for the sample mean, whose expectation is 1 when there is no interference. For each data point, we first average 250 pulses using the method indicated in the title and then average four such groups arithmetically for a total of 1000 pulses. Using a smaller number of pulses makes the memory requirement less stringent and the trimmed mean more efficient. The ionosphere signal is largely characterized by smooth temporal and spatial varia-





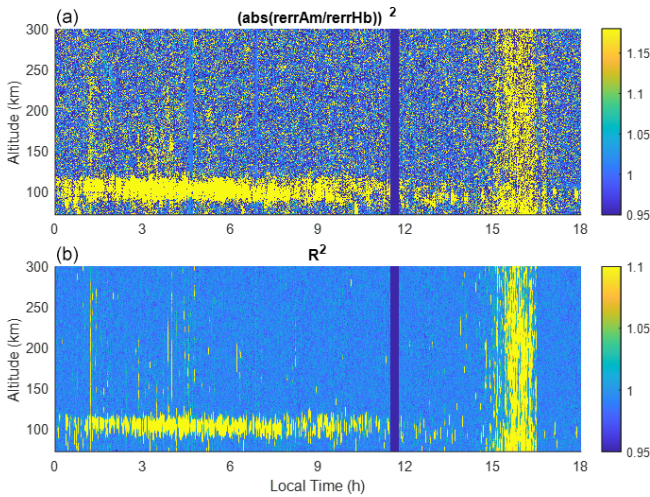
**Figure 1.** Range–time–intensity (RTI) plots of incoherent scatter total power returns on 11–12 September 2014. The first four panels, starting from the top, are the power return of the sample mean, the trimmed mean at the 95 % level, the trimmed mean based on the geometric deviation, and a hybrid method, respectively. The last panel is the normalized standard deviation.

tions during the daytime and by thin horizontal layers, known as sporadic E's, around 100 km at nighttime. The study of sporadic E layers and the associated dynamics have attracted much attention and are active areas of research (e.g., Mathews, 1998; Larsen et al., 2007; Wang et al., 2022; Kunduri et al., 2023). Two types of interference seen in Fig. 1 are represented in boxes A and B in the top panel. Box A is likely another radar operating at the same inter-pulse period (IPP) as that of the Arecibo ISR or is an internal system problem. Vertical lines in box B and other similar vertical lines that are confined to  $\sim 90$ –120 km are meteoric echoes. The altitude extension of meteor echoes is because fast-moving meteor heads cannot be decoded by the matched filter. They do not extend beyond 120 km in altitude in our case because meteor echoes are detected below about 115 km (Zhou and Kelley, 1997). The normalized standard deviation  $R(A_N)$  is displayed in the bottom panel in Fig. 1.

The top panel in Fig. 1 shows the result of arithmetically averaging 1000 pulses (i.e., the sample mean). All types of interference show up prominently, as the method does not fil-

ter out any contamination. The trimmed mean (second panel) cleans up the first part of the heavy contamination in box A but is not effective against the second part, most likely because more than 5 % of the pulses were contaminated.  $T_{\text{GEO4}}$  and the hybrid method largely filters out the contamination in box A and reveal the underlying sporadic layer despite the heavy contamination. Although  $T_{\text{GEO4}}$  appears to handle all the contamination as well as the hybrid method does, it is slightly inferior to the latter in reducing statistical error, as seen in the later part of this section. The only residue contamination not filtered out is around 22:30 LT. None of the methods is effective in removing it completely, and all three robust estimators appear to perform the same. As the total power of the interference is relatively low, interference may permeate most of the pulses, making it very difficult to remove it from each pulse. For this type of interference, one way is to find the mean at non-ionosphere heights and subtract it from the entire profile. Noise samples are available at Arecibo. Background noise is not subtracted here to focus on the effect of robust estimators in this study. The trimmed



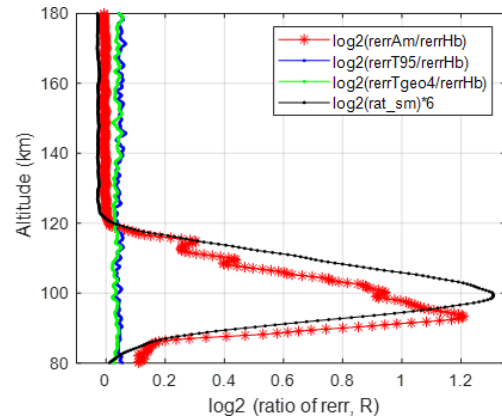


**Figure 2.** (a) The square of the relative error in the sample mean method normalized to that of the hybrid method. (b) The normalized variance. The yellow color in (a) indicates that the sample mean has a larger error than the hybrid method.

mean,  $T_{GEO4}$ , and hybrid methods are all effective at removing meteor interference, which typically does not last more than 50 ms at Arecibo, i.e., 5 pulses (Zhou and Kelly, 1997).

Other than the most obvious interference highlighted in boxes A and B, no other contamination appears to be obvious. The  $R$  value in the region indicated by box C has elevated values, indicating likely contamination. Yet there appears to be little difference between the sample mean result in the top panel of Fig. 1 and the results from the robust estimators. One effect of the interference is that it increases the statistical error, which is more difficult to see from the RTI plot. To estimate the statistical error, we use the difference in the power minus the average power of the surrounding 15 points in height and 5 points in time as a proxy for the error. The square ratio of the sample mean error to the error in the hybrid method is displayed in Fig. 2a. The corresponding  $R(A_N)$  is displayed in Fig. 2b. Larger statistical error from the sample mean in the region indicated by box C in Fig. 1 is quite evident. Although  $R(A_N)$  is not linearly related to the error, elevated  $R(A_N)$  is a robust indicator of contamination. This is also evidenced from 01:00 to 03:00 LT in Fig. 2, where sporadic elevations of  $R(A_N)$  and statistical errors are seen to be correlated.

An estimator needs to be efficient when there is no interference. Figure 3 shows the ratio of the sample mean and  $T_{95}$  errors to the hybrid error as well as the corresponding standard deviation  $R(A_N)$  averaged between 07:00 to 13:00 LT, during which period contamination is minimal above 120 km (as seen in Fig. 2). The error in the hybrid estimator is virtually the same as that in the arithmetic average. The error in the  $T_{95}$  estimator is 1.036 times the error in the hybrid estimator, which is in good agreement with the simulated value of  $\sqrt{1.09/1.018} = 1.035$ . Similarly, the error in  $T_{GEO4}$  is



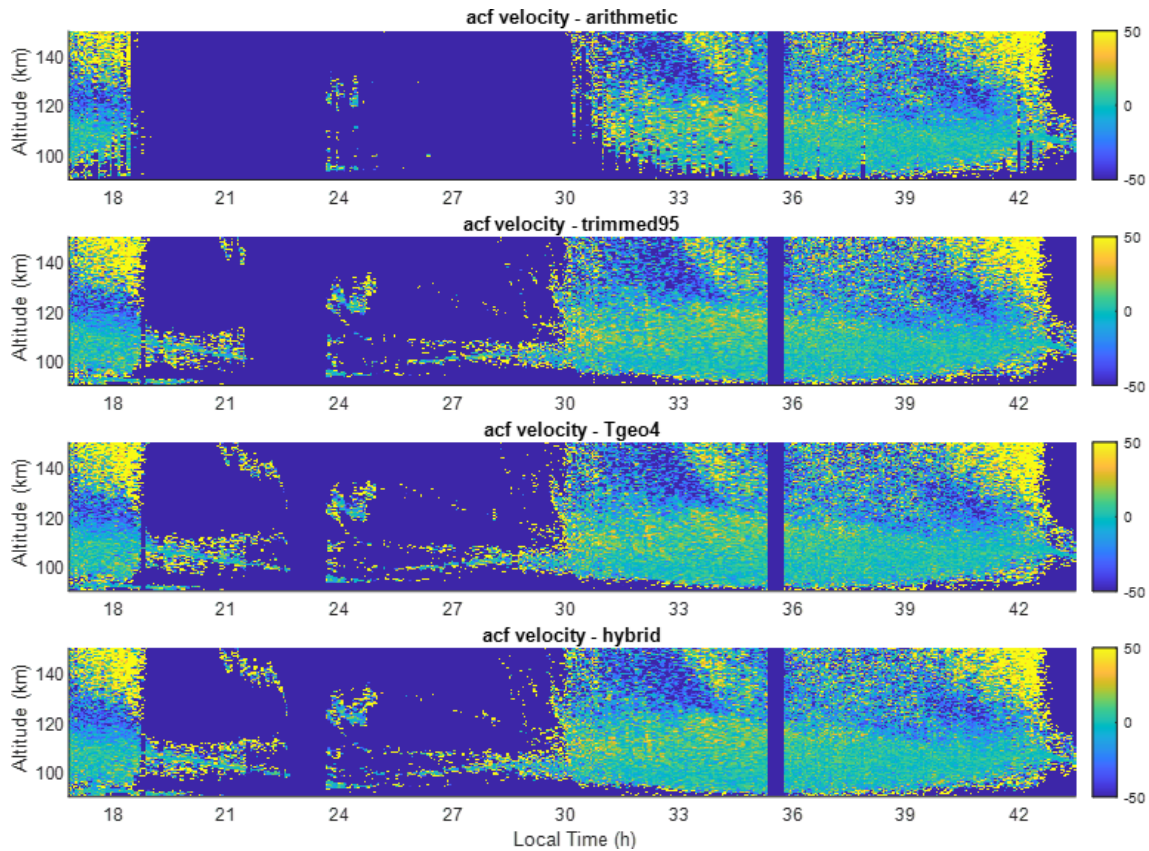
**Figure 3.** Mean relative errors (in base-2 logarithms) of the sample mean, trimmed mean, and  $T_{GEO4}$  normalized to that of the hybrid method (red and blue lines, respectively). The black line is 6 times the logarithm (base 2) of the mean  $R$ . The time duration averaged for all the lines in the figure is from 07:00 to 13:00 LT on 12 September 2014.

slightly smaller than that in  $T_{95}$ , which is also in good agreement with the simulation results shown in Table 1. The mean  $R(A_N)$  correlates with the elevated error in the region of 90–120 km. We also note that the mean  $R(A_N)$  above 120 km is 0.997, which is slightly below the expected value of 1. Although the deviation is small, it is statistically significant. This may be caused by the bias in the receiving channels or the finite dynamic range of the analog-to-digital converters.

### 3.2 Power spectrum processing and Doppler velocity comparisons

The power spectral density (PSD) of an ISR is obtained by transmitting a coded long pulse (CLP), 440  $\mu$ s in our case. The baud length is 2  $\mu$ s, making the bit number of the pulse 220. The inter-pulse period is 10 ms as in the Barker data. The bit sequence is random for each transmitted pulse. The PSD is obtained by the Fourier transform of the data multiplied by the complex conjugate of the code. The characteristics of the CLP are discussed by Sulzer (1986). The averaging of the PSD at each frequency component is identical to that of the total power in the above section, which can be viewed as the center frequency component.

Figure 4 shows the Doppler velocity derived from the four estimators using the phase of the auto-correlation function. The vertical ion drift in the altitude range of 90–150 km is typically less than 50  $\text{ms}^{-1}$  above Arecibo. Below 120 km, the plasma drift is the same as the neutral wind because of the complete coupling between ions and neutral molecules. During the daytime, there are sufficient signals above 95 km to obtain continuous spatial and temporal velocities. During the nighttime, it is only possible to obtain velocities within thin ionization layers. While ion velocity with fine height and time resolutions is of great geophysical interest (e.g., Zhou



**Figure 4.** Vertical ion velocities obtained using the four estimators. The estimators, from top to bottom, are the sample mean, trimmed mean (at 95 %),  $T_{GEO4}$ , and hybrid mean, respectively.

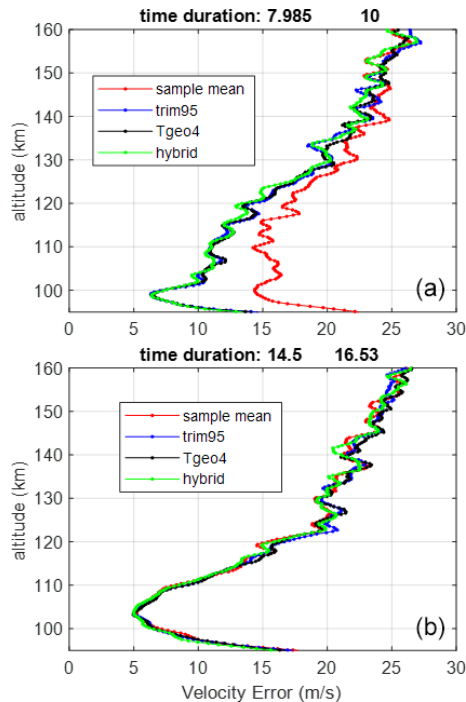
et al., 1997; Hysell et al., 2014), our focus here is to study the relative accuracy of the velocities obtained from different estimators.

Comparisons of the velocity results largely follow those of the total power. The sample mean fails in boxes A and B. Additionally, during the sunrise hours when the ionospheric signal is low and the meteoric interference is strong, the sample mean can only yield valid velocities occasionally while the robust estimators can obtain the velocities continuously in altitude and time. As in the total power estimation, the trimmed mean does not yield valid results in the second part of box A from 21:30 to 22:30 LT, while the hybrid and  $T_{GEO4}$  methods appear not to be affected by the interference very much.

To compare the statistical fluctuations, we use the altitudinal difference in the velocity divided by the square root of 2 as a proxy for velocity error. Figure 5 shows the altitude variation in the velocity error during 08:00–10:00 LT as well as 14:30–16:30 LT on 12 September. All the robust estimators have essentially the same error at each altitude, while the sample mean has a much larger error around 100 km. The error in the sample mean converges to those of the robust estimators above 145 km. The diminishing error difference in the sample mean with increasing altitude is due to the long pulse

length (440  $\mu$ s) used. A characteristic of the CLP pulse is that the interference at one altitude is uniformly spread across the entire bandwidth randomly at other altitudes. A meteor echo at 100 km increases the spectral power fluctuations with diminishing strength up to 166 km. Meteoric influx peaks at 06:00 LT and varies strongly with the local time. The daily variation in meteoric flux is quantitatively analyzed by Zhou et al. (1995) and Li and Zhou (2019). It can also be qualitatively seen in Fig. 2b. The larger error in the sample mean during 08:00–10:00 LT is a reflection of the strong meteoric flux. Although the afternoon period suffers from meteoric interference and radio contamination, as seen in Fig. 2, both of them are weak. Statistical averaging of 6000 pulses is able to even out the spectral power fluctuation to such a degree that all the estimators produce the same velocity. For spectral processing, the most important factor is the total amount of noise power, while the percentage of pulses contaminated is often more important in total power processing.

Overall, we see that the  $T_{GEO4}$  and hybrid estimators accurately and consistently improve velocity and total power measurements over the sample mean, which are important for studying the E-region dynamics and composition. The availability of nighttime velocities will help reduce the large



**Figure 5.** Doppler velocity errors for the sample mean, trimmed mean,  $T_{\text{GEO4}}$ , and hybrid method on 12 September 2014 for 08:00–10:00 LT (a) and 14:30–16:30 LT (b).

error in the measurement of atmospheric tides in the E region (Zhou et al., 1997; Gong et al., 2013). Accurate measurement of the power spectrum and total power will facilitate all E-region studies, especially those concerning the climatology and dynamics of sporadic E and intermediate layers (Zhou et al., 2005; Hysell et al., 2009; Raizada et al., 2018; Gong et al., 2021). Of particular importance are the vertical wind and ion composition of the E region, which have not been studied much due to a lack of quality data.

#### 4 Summary and conclusion

We have discussed several robust estimators to compute the variance of a normally distributed random variable,  $X$ , to deal with interference. This variance is the same as the mean of the power variable,  $X^2$ . The effectiveness of an estimator is described by the normalized standard deviation,  $R$ . We derive the theoretical  $R$  values for the median, geometric mean, and trimmed mean of gamma distributions, which result from averaging the power random variables. We discuss and compare another four estimators through simulations for various interference scenarios. Robust estimators found in the literature are typically rank-based (e.g., the median, trimmed mean, and median absolute deviation). We have used the geometric mean and geometric deviation as two basic parameters in assessing the likelihood of a data point being contaminated. The methods based on the geometric mean have

two advantages over the rank-based ones: they are less susceptible to the large uncertainties in the tail part of the distributions and they are computationally more efficient. For the interference model used, the  $T_{\text{GEO4}}$  estimator, which is based on the geometric mean, is particularly effective as a stand-alone estimator when there is no initial average. Another effective estimator based on the geometric mean is the weighted mean. The  $R$  value of the sample mean can be used to assess whether the process conforms to the expected distribution. This knowledge allows us to combine the sample mean with other robust estimators to mitigate contamination and achieve statistical accuracy.

We apply three robust estimators to incoherent scatter power and velocity processing, along with the traditional sample mean estimator. We show that the performance of estimators with real data agrees well with simulations. In the total power processing, the trimmed mean performs mostly well except when the contamination is very heavy. The  $T_{\text{GEO4}}$  estimator performs almost as well as the hybrid method in mitigating interference. The hybrid method performs the best at mitigating interference as well as at reducing statistical errors. For Doppler velocity processing, the same conclusion can be drawn in cases of frequent interference. When the interference is weak, all the robust estimators appear to perform well. For the Arecibo ISR data, the sample mean has larger statistical errors even for data that may not appear to contain obvious interference. This highlights the need for robust estimation to process or reprocess decades of E-region data taken at Arecibo. The hybrid estimator is most advantageous under all circumstances. This conclusion is likely applicable to other incoherent scatter radars as well. While the interference characteristics differ at each radar site, the study provides a foundation to optimize robust estimation, which is an essential step in many data processing applications.

*Data availability.* The raw data are available upon request from the Texas Advanced Computing Center (<https://tacc.utexas.edu/research/tacc-research/arecibo-observatory/>, last access: August 2023). The processed and the raw data can be obtained from the corresponding author.

*Author contributions.* QZ conceptualized the framework, formulated the theory, and drafted and finalized the paper. YL did the data analysis and contributed to the data visualization. YG contributed to the conceptualization and verification of the theory and edited the paper.

*Competing interests.* The contact author has declared that none of the authors has any competing interests.

*Disclaimer.* Publisher's note: Copernicus Publications remains neutral with regard to jurisdictional claims made in the text, published maps, institutional affiliations, or any other geographical representation in this paper. While Copernicus Publications makes every effort to include appropriate place names, the final responsibility lies with the authors.

*Acknowledgement.* We thank the former staff members of the Arecibo Observatory for collecting the data.

*Financial support.* This research has been supported by the US National Science Foundation, Directorate for Geosciences (grant no. 2152109).

*Review statement.* This paper was edited by Meng Gao and reviewed by two anonymous referees.

## References

- Chau, J. L., Woodman, R. F., and Galindo, F.: Sporadic meteor source as observed by the Jicamarca high-power large-aperture VHF radar, *Icarus*, 188, 162–174, <https://doi.org/10.1016/j.icarus.2006.11.006>, 2007.
- Evans, J. V.: Theory and practice of ionosphere study by Thomson scatter radar, *P. IEEE*, 57, 496–530, 1969.
- Gong, Y., Zhou, Q., and Zhang, S.: Atmospheric tides in the low latitude E- and F-region and their response to a sudden stratospheric warming event in January 2010, *J. Geophys. Res.-Space*, 118, 7913–7927, <https://doi.org/10.1002/2013JA019248>, 2013.
- Gong, Y., Lv, X., Zhang, S., Zhou, Q., and Ma, Z.: Climatology and seasonal variation of the thermospheric tides and their response to solar activities over Arecibo, *J. Atmos. Sol.-Terr. Phys.*, 215, 105592, <https://doi.org/10.1016/j.jastp.2021.105592>, 2021.
- Huber, P. J. and Ronchetti, E. M.: *Robust Statistics*, 2nd edn., John Wiley and Sons, ISBN 978-0-470-12990-6, 2009.
- Hysell, D. L., Nossa, E., Larsen, M. F., Munro, J., Sulzer, M. P., and Gonzalez, S. A.: Sporadic E layer observations over Arecibo using coherent and incoherent scatter radar: Assessing dynamic stability in the lower thermosphere, *J. Geophys. Res.*, 114, A12303, <https://doi.org/10.1029/2009JA014403>, 2009.
- Hysell, D. L., Larsen, M. F., and Sulzer M. P.: High time and height resolution neutral wind profile measurements across the mesosphere/lower thermosphere region using the Arecibo incoherent scatter radar, *J. Geophys. Res.-Space*, 119, 2345–2358, <https://doi.org/10.1002/2013JA019621>, 2014.
- Isham, B., Tepley, C., Sulzer, M., Zhou, Q., Kelley, M., Friedman, J., and Gonzalez, S.: Ionospheric observations at the Arecibo Observatory: Examples obtained using new capabilities, *J. Geophys. Res.*, 105, 18609, <https://doi.org/10.1029/1999JA900315>, 2000.
- Kunduri, B. S. R., Erickson, P. J., Baker, J. B. H., Ruohoniemi, J. M., Galkin, I. A., and Sterne, K. T.: Dynamics of mid-latitude sporadic-E and its impact on HF propagation in the North American sector, *J. Geophys. Res.-Space*, 128, e2023JA031455, <https://doi.org/10.1029/2023JA031455>, 2023.
- Larsen, M. F., Hysell, D. L., Zhou, Q. H., Smith, S. M., Friedman, J., and Bishop R. L.: Imaging coherent scatter radar, incoherent scatter radar, and optical observations of quasiperiodic structures associated with sporadic E layers, *J. Geophys. Res.*, 112, A06321, <https://doi.org/10.1029/2006JA012051>, 2007.
- Li, Y. and Zhou Q.: Characteristics of micrometeors observed by the Arecibo 430 MHz incoherent scatter radar, *Mon. Not. R. Astron. Soc.*, 486, 3517–3523, <https://doi.org/10.1093/mnras/stz1073>, 2019.
- Li, Y., Galindo, F., Urbina, J., Zhou, Q., and Huang T.-Y.: A Machine Learning Algorithm to Detect and Analyze Meteor Echoes Observed by the Jicamarca Radar, *Remote Sens.-Basel*, 15, 4051, <https://doi.org/10.3390/rs15164051>, 2023.
- Mathews, J. D.: Sporadic E: current views and recent progress, *J. Atmos. Sol.-Terr. Phys.*, 60, 413–435, 1998.
- Raizada, S., Brum, C. G., Mathews, J. D., Gonzalez, C., and Franco, E.: Characteristics of nighttime E-region over Arecibo: Dependence on solar flux and geomagnetic variations, *Adv. Space Res.*, 61, 1850–1857, <https://doi.org/10.1016/j.asr.2017.07.006>, 2018.
- Resnick, S. I.: *Heavy-Tail Phenomena: Probabilistic and Statistical Modeling*, Springer, New York, ISBN-10 0-387-24272-4, 2007.
- Rousseeuw, P. J. and Croux, C.: Alternatives to the Median Absolute Deviation, *J. Am. Stat. Assoc.*, 88, 1273–1283, 1993.
- Siegrist, K.: *Probability, Mathematical Statistics, Stochastic Processes*, LibreTexts, [https://stats.libretexts.org/Bookshelves/Probability\\_Theory/Probability\\_Mathematical\\_Statistics\\_and\\_Stochastic\\_Processes\\_\(Siegrist\)](https://stats.libretexts.org/Bookshelves/Probability_Theory/Probability_Mathematical_Statistics_and_Stochastic_Processes_(Siegrist)) (last access: 28 April 2024), 2022.
- Stigler, S. M.: The asymptotic distribution of the trimmed mean, *Ann. Stat.*, 1, 472–477, 1973.
- Sulzer, M. P.: A radar technique for high range resolution incoherent scatter autocorrelation function measurements utilizing the full average power of klystron radars, *Radio Sci.*, 21, 1033–1040, 1986.
- Taleb, N. N.: *Statistical Consequences of Fat Tails: Real World Preasymptotics, Epistemology, and Applications*, STEM Academic Press, ISBN 978-1-5445-0805-4, 2022.
- Wang, Y., Themens, D. R., Wang, C., Ma, Y.-Z., Reimer, A., Varney, R., Gilies, R., Xing, Z.-Y., Zhang, Q.-H., and Jayachandran, P. T.: Simultaneous observations of a polar cap Sporadic-E layer by twin incoherent scatter radars at Resolute, *J. Geophys. Res.-Space*, 127, e2022JA030366, <https://doi.org/10.1029/2022JA030366>, 2022.
- Wikipedia: Gamma distribution, [https://en.wikipedia.org/wiki/Gamma\\_distribution](https://en.wikipedia.org/wiki/Gamma_distribution), last access: 28 April 2024a.
- Wikipedia: Gamma function, [https://en.wikipedia.org/wiki/Gamma\\_function](https://en.wikipedia.org/wiki/Gamma_function), last access: 28 April 2024b.
- Wilcox, R.: *Introduction to Robust Estimation and Hypothesis Testing*, 4th edn., Elsevier, Amsterdam, ISBN 978-0-12-804733-0, 2017.
- Zhou, Q. and Kelley, K. C.: Meteor observation by the Arecibo 430 MHz ISR II. results from time-resolved observations, *J. Atmos. Sol.-Terr. Phys.*, 59, 739–752, [https://doi.org/10.1016/S1364-6826\(96\)00103-4](https://doi.org/10.1016/S1364-6826(96)00103-4), 1997.
- Zhou, Q., Tepley, C. A., and Sulzer, M. P.: Meteor observations by the Arecibo 430 MHz incoherent scatter radar – Results from time-integrated observations, *J. Atmos. Terr. Phys.*, 57, 421–432, [https://doi.org/10.1016/0021-9169\(94\)E0011-B](https://doi.org/10.1016/0021-9169(94)E0011-B), 1995.

Zhou, Q., Sulzer, M. P., and Tepley, C. A.: An analysis of tidal and planetary waves in the neutral winds and temperature observed at the E-region, *J. Geophys. Res.*, 102, 11491–11505, 1997.

Zhou, Q., Friedman, J., Raizada, S., Tepley, C. A., and Morton, Y. T.: Morphology of nighttime ion, potassium and sodium layers in the meteor zone above Arecibo, *J. Atmos. Sol.-Terr. Phys.*, 67, 1245–1257, <https://doi.org/10.1016/j.jastp.2005.06.013>, 2005.

Zhou, Q. H., Zhou, Q. N., and Mathews, J. D.: Arithmetic average, geometric average and ranking – Application to incoherent scatter radar data processing, *Radio Sci.*, 34, 1227–1237, 1999.