

Fitting Wind Speed to a Probability Distribution Using Minimum Distance Estimation Technique

Otieno Okumu Kevin^{*}, John Matuya, Muthiga Nganga

Department of Mathematics and Physical Sciences, Maasai Mara University, Narok, Kenya

Email address:

Kevinotieno15@gmail.com (O. O. Kevin), jmatuya@mmarau.ac.ke (J. Matuya), muthiga@mmarau.ac.ke (M. Nganga)

^{*}Corresponding author

To cite this article:

Otieno Okumu Kevin, John Matuya, Muthiga Nganga. Fitting Wind Speed to a Probability Distribution Using Minimum Distance Estimation Technique. *American Journal of Theoretical and Applied Statistics*. Vol. 10, No. 6, 2021, pp. 226-232.

doi: 10.11648/j.ajtas.20211006.11

Received: March 10, 2021; **Accepted:** March 22, 2021; **Published:** November 10, 2021

Abstract: From the past studies, we realized that minimum distance estimation technique is not commonly used for fitting wind speed data to a distribution yet it is believed to be the best alternative for Maximum Likelihood Estimation (MLE) method which is known to give good estimates than Least Square Estimates (LSE) and Method of Moments (MOM). To achieve this, the study aims at fitting data to a probability distribution using minimum distance estimation techniques to find the best distribution. The study uses wind speed data from five sites in Narok county namely; Irbaan primary, Imortott primary, Mara conservancy, Oldrkesi and Maasai Mara University. The best wind speed models were examined using the Cullen and Frey graph and a suitability test on the models done using Kolmogorov-Smirnov statistical test of goodness of fit. The wind speed data are fitted to the recommended distributions using minimum distance estimation techniques. The best distribution was identified using Akaike's Information Criterion (AIC) and Bayesian Information criterion (BIC). From the distribution comparison for the two and three parameter distributions, gamma is the best in all cases. Gamma with three parameter distribution gives lower AIC and BIC values and model comparison test showing that gamma 3-parameter is the better than gamma with 2-parameters. The study concluded that gamma distribution with three parameters is the best distribution for fitting wind speed data with the three parameters given as; threshold parameter of 0.1174, shape parameter of 1.8646 and scale parameter of 0.9937.

Keywords: Minimum Distance Estimation (MDE), Akaike's Information Criterion (AIC), Bayesian Information Criterion (BIC), Distribution

1. Introduction

Many researchers have fitted several distributions for predicting wind speeds. Some of the fitted distributions are Weibull, gamma, log-normal, Rayleigh distribution, hazard Weibull, and Erlang among others [1, 7, 16]. In Kenya, a study on the wind speed distribution with two parameters in studying wind regime analysis and reserve estimation and in analysis of wind behavior in Juja site respectively was conducted using Weibull distribution with two parameters [2, 15]. A study on Weibull distribution with two parameters found that the Weibull parameters which are scale and shape parameters are used to show how the site is windy and the variability of the site in terms of peakedness [16]. This wind speed

distributions are examined both for two and three parameters. The two parameters are scale and shape parameters while the three parameters distribution has scale, shape and threshold parameter. From the several distributions, the most commonly used distributions from the past studies are Weibull distribution, gamma distribution and log-normal distribution for both two parameters and three parameters distributions.

The researchers who studied different wind speed distributions used different fitting techniques for estimating the model parameters. Some of the common techniques applied on the analysis process are maximum likelihood estimation, method of moment estimation,

least square estimation method, mean wind speed and standard deviation (point estimation) method [21]. A study on comparison of maximum likelihood estimation method, least square estimation method and method of moment estimation method and found that maximum likelihood estimation method gives goods estimates compared to moment method and least square estimations method [18].

From the past related studies, a promising alternative to maximum likelihood estimation techniques is the minimum distance estimation technique because it is considered less sensitive to the problematic assumption of maximum likelihood estimation [13]. Thus, it is referred to a robust estimation technique implying that it attempts to protect against minor deviations from the underlying assumptions. The concept of minimum distance estimate is that better estimates will be obtained by fitting a distribution to a sample data. From the assessment of related studies, it was found that researchers are not applying the minimum distance technique in fitting the single 2-parameter or 3-parameter distributions.

1.1. Two Parameter Distributions

1.1.1. Weibull Distribution

Researchers who has used Weibull distribution to analyze wind speed concluded that the Weibull distribution function is the best in estimating the parameters of wind speed. The Weibull distribution model applied by the researcher is given by [1, 6, 20];

$$f(u) = \frac{b}{p} \left(\frac{u}{p} \right)^{b-1} \exp \left[- \left(\frac{u}{p} \right)^b \right], (b, u > 0, p > 1), \quad (1)$$

Where:

$f(u)$ is the probability of observing wind speed,

u is the wind speed,

b is the shape factor (parameter) which has no unit but range from 1.5 to 3.0 for most wind conditions,

p is the value in the unit of wind speed called the Weibull scale parameter in m/s.

1.1.2. Lognormal Distribution

Wind speed analysis is very wide and one of the statistical distributions used in examining the wind data is the log-normal statistical model with parameters v and k [3, 22]. The log-normal density function with the two parameters is given by:

$$f(p) = \frac{1}{k\sqrt{2p\pi}} \exp \left[- \frac{(\ln p - v)^2}{2k^2} \right], \quad (2)$$

Where:

p is the log-normal random variable,

$\ln(p)$ is the normal random variable,

v is the mean for a normal random variable,

k is the standard deviation for the normal random variable.

1.1.3. Gamma Distribution

The probability density function of gamma random variable y in combination with two parameters z and q representing the shape and scale parameters respectively is given by the past studies [10, 12].

$$f(y, z, q) = \frac{y^{z-1} \exp \left(- \frac{y}{q} \right)}{\Gamma(z) q^z}, (z, y, q, > 0) \quad (3)$$

Where:

$$\Gamma(v) = \int_0^\infty y^{v-1} \exp^{-y} dy, v > 0 \quad (4)$$

And:

z is the shape parameter,

q is the scale parameter,

y are the random variables (wind speed).

1.2. Three Parameter Distributions

1.2.1. Weibull Distribution with 3 Parameters

The Weibull statistical distribution with three parameters is given by [1, 4, 8];

$$f(u) = \left(\frac{b}{p} \right) \left(\frac{u-w}{p} \right)^{b-1} \exp \left[- \left(\frac{u-w}{p} \right)^b \right] (b, u > 0, p > 1), \quad (5)$$

Where:

u is the wind speed,

b is the shape parameter,

p is the scale parameter measured in m/s,

w is the thresh-hold parameter.

1.2.2. Lognormal Distribution with 3 Parameters

This distribution has three parameters namely scale parameter, shape parameter and thresh-hold parameter also known as location parameter. The probability density function and the cumulative density function are given by the below equations [5, 10],

$$f(p) = \frac{1}{(p-y) k\sqrt{2\pi}} \exp \left[- \left(\frac{\ln(p-y)-v}{2k} \right)^2 \right] \quad (6)$$

Where:

$v > 0$ is the scale parameter,

$k > 0$ is the shape parameter,

y is the thresh-hold parameter, also referred to the location parameter,

$p \geq$ is the wind speed.

1.2.3. Gamma Distribution with 3 Parameters

From the past studies the gamma function is given as follows [11, 14, 20];

$$f(y, z, q, t) = \frac{y t^{z-1}}{\Gamma(z/t) q^z} \exp \left[- \left(\frac{y}{q} \right)^t \right] (z, y, q, t > 0), \quad (7)$$

Where:

q is the scale parameter,

z, is shape parameters,

t is the thresh-hold parameter,

2. Data and Methods

The data comprises of 63778 hourly wind speed observations. The minimum speed in the data is 0.12 m/s and the maximum speed is 5.35 m/s with the mean speed of 1.9658 m/s and the estimated standard deviation of 1.2407 m/s. The estimated kurtosis and skewness are 2.8094 and 0.8433.

2.1. Cullen and Frey Graph

This graph helps in understanding the best possible distribution or distributions that the data is fitting best.

From Figure 1, we can observe that the plot is can be estimated at around a kurtosis of 2.8 and square of skewness of around 0.7 (skewness = 0.8). With a skewness of 0.8 and kurtosis of 2.8 we can conclude that the normal distribution cannot fit the data best since normal distribution requires that kurtosis = 3 and skewness = 0. The uniform distribution is not also the best distribution for fitting this data since the observed difference between the scatter plot of kurtosis and square of skewness and that of uniform distribution is not that close (for a uniform distribution needs a kurtosis value of 1.8 and skewness value of 0).

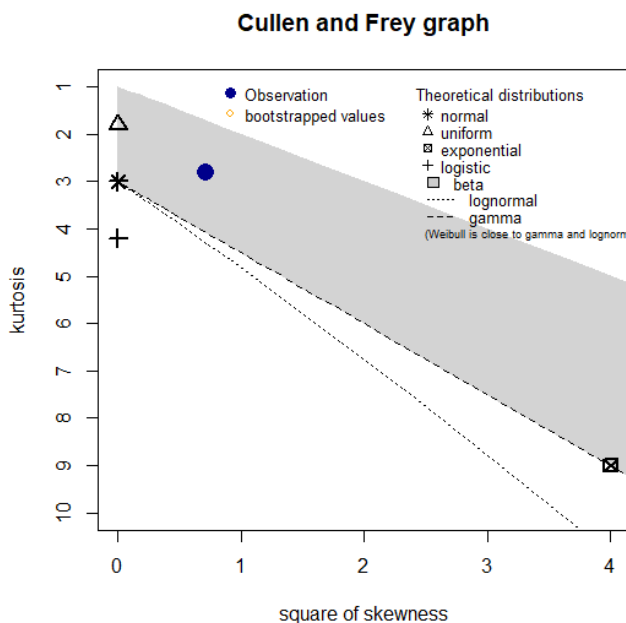


Figure 1. Cullen and Frey graph.

For the logistic distribution, we can say that it is also not the best for the data since logistic distribution always have a kurtosis of 4.2 and skewness value = 0 and from the graph it can be observed that the logistic plot is not closer to the data plot. For the exponential distribution we can observe that its point is far away from the data point, this is because exponential distribution is expected to have a kurtosis of 9

and skewness of 2 compared to the data point skewness of 0.7 therefore exponential distribution is not the best for the data. From the graph, it can be seen that beta distribution can fit the data but this distribution cannot be applied to the data since beta distribution is a family of continuous probability distributions defined on the interval of [0, 1] which is not the case with the collected data for this research. From the graph, log-normal and gamma distributions can fit the data best because they appear to be close to the data points and well distributed. Weibull distribution is also another good distribution for fitting the data since from the graph it is said that Weibull is close to gamma and log-normal.

2.2. Data Description After Subtracting the Threshold Value

The value of the threshold tries to plot the data to a straight line when subtracted from the original data. The threshold values take the same unit as the units for the examined data set and is used as a transformation to reduce biasedness. The threshold value used for the three-parameter analysis is 0.1174. After subtracting the threshold value, the statistic for the data is as shown in table 1.

Table 1. Summary statistics after subtracting the threshold value.

Min value	0.0026
Max value	5.2326
Mean	1.8484
Median	1.5026
Estimated std	1.24065
Estimated kurtosis	2.809401
Estimated skewness	0.8433485

2.3. Minimum Distance Estimation Technique

Minimum distance method reduces the computational complexity since it omits the jacobian element which is usually present in the likelihood function [17].

2.3.1 Concept of Minimum Distance Estimation

If given a statistics $\hat{\beta}$ with the distribution $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \theta)$ We are also given a distance $k(\cdot, \cdot)$ which is continuously differentiable. We assume that there is a unique parameter α such that

$$k(\beta, \alpha) = 0 \quad (8)$$

We have that k is $L \times 1$, α is $M \times 1$ and β is $L \geq M$. The minimum distance solves the following

$$\hat{\alpha} = \arg \min (k(\hat{\beta}, \alpha)' \hat{C} k(\hat{\beta}, \alpha)) \quad (9)$$

Where \hat{C} converges in probability to a $L \times L$ non-random, positive definite and symmetric matrix C . The first order condition for the minimum distance method is given as

$$(\partial k(\hat{\beta}, \hat{\alpha}) / \partial \alpha)' \hat{C} k(\hat{\beta}, \hat{\alpha}) = 0 \quad (10)$$

Therefore, we can write the minimum distance estimator as satisfying the following equation

$$\bar{D}k(\hat{\beta}, \hat{\alpha}) = 0, \bar{D} = (\partial k(\hat{\beta}, \hat{\alpha}) / \partial \alpha)' \hat{C} \quad (11)$$

The method of minimum distance estimation technique

depends on the test statistics of Anderson-Darling (AD) test [9]. The expression for the two methods based on minimum distance estimation is formulated as follows

$$AD = A_n^2 = -n - \frac{1}{n} \sum_{i=1}^n (2i-1) [\log y_{i:n} + \log(1 - y_{(n+1-i)})] \quad (12)$$

For the two parameter estimates for Weibull, lognormal and gamma distributions, Anderson-Darling was applied to assist in estimating the two-parameters using the minimum distance method.

2.3.2. Method of Anderson-Darling Estimation

A good approach of minimum distance estimator is based on the application of Anderson-darling statistics and is defined as Anderson-Darling estimator (ADE).

Minimum Distance Estimation was developed as an alternative to statistical test to be used significantly to examine sample distribution departure from normality [23]. By applying the Anderson-Darling test statistics, we can obtain the Anderson-Darling estimates \hat{m}_{ADE} and \hat{k}_{ADE} representing the scale and shape parameter estimates respectively for the three distributions from the following equation.

$$A(m, k) = -n - \frac{1}{n} \sum_{i=1}^n (2i-1) \{\log F(y_{i:n}|m, k) + \log V(y_{(n+1-i)}|m, k)\} \quad (13)$$

The estimates \hat{m}_{ADE} and \hat{k}_{ADE} are obtained by minimizing equation (12) with respect to m and k. Similarly, these estimates can be obtained from the solution of the following nonlinear equations [19].

$$\begin{aligned} \sum_{i=1}^n (2i-1) \left[\frac{\Delta_1(y_{i:n}|m, k)}{F(y_{i:n}|m, k)} - \frac{\Delta_1(y_{(n+1-i)}|m, k)}{V(y_{(n+1-i)}|m, k)} \right] &= 0 \\ \sum_{i=1}^n (2i-1) \left[\frac{\Delta_2(y_{i:n}|m, k)}{F(y_{i:n}|m, k)} - \frac{\Delta_2(y_{(n+1-i)}|m, k)}{V(y_{(n+1-i)}|m, k)} \right] &= 0 \end{aligned} \quad (14)$$

Where, $\Delta_1(\cdot|m, k)$ and $\Delta_2(\cdot|m, k)$ are given in equation (14) below

$$\begin{aligned} \Delta_1(y_{i:n}|m, k) &= \frac{\partial}{\partial m} F(y_{i:n}|m, k) = \frac{e^{ky_i} - 1}{(e^{ky_i} - 1 + m)^2} \\ \Delta_2(y_{i:n}|m, k) &= \frac{\partial}{\partial k} F(y_{i:n}|m, k) = \frac{ky_i e^{ky_i}}{(e^{ky_i} - 1 + m)^2} \end{aligned} \quad (15)$$

After getting the efficient and precise threshold values based on the AIC and BIC values. The same technique of Anderson Darling estimation is applied for the estimating the two-parameter distributions namely: scale and shape of Weibull, lognormal and gamma.

2.4. Kolmogorov-Smirnov Test

This is a two-sample test with the advantage that it does not depend mostly on the underlying cumulative distribution function being tested and also applies only to continuous distributions which in this case is applicable since we are only investigating the continuous statistical distributions [7]. It is calculated as;

$$D^* = \max(|F_1(t) - F_2(t)|) \quad (16)$$

Where:

$F_1(t)$ is the proportion of t1 values less than or equal to t

$F_2(t)$ is the proportion of t2 values less than or equal to t

H_0 : The data follows specified distribution

H_1 : The data do not follow specified distribution

The smaller the test statistic the better the fit.

2.5. Akaike's Information Criterion (AIC)

The Akaike's Information Criterion is calculated as;

$$AIC = -2\log L(p) + 2w \quad (17)$$

Where $\log L(P)$ defines the value of the maximized log-likelihood objective function for a model with w parameters. A smaller AIC value represents a better fit.

2.6. Bayesian Information Criterion (BIC)

The Bayesian Information Criterion is calculated as;

$$BIC = -2\log L(p) + w \log M \quad (18)$$

Where $\log L(P)$ represents the values of the maximized log-likelihood objective function for a model with w parameters fit to M data points. A smaller Bayesian Information Criterion value indicates a better fit (best model for fitting the data).

3. Results

3.1. Analysis for Two Parameter Distributions

By using the distance technique method, the estimated parameters for the three distributions in examination is given as in the table 2.

Table 2. Parameters.

Distribution	Parameter	Estimate
Weibull	Shape	1.502943
	Scale	2.172747
Gamma	Shape	2.107526
	Scale	1.062046
Log-normal	Shape	0.490416
	Scale	0.68618

Also, under this technique there is need to identify the best distribution for fitting the data under study. The table 3, shows the goodness test of fit criteria and goodness test of fit statistics. The goodness of test statistic is used to test if the distribution can fit the data under study while the goodness of fit criteria is used to understand the best distribution for the data.

Table 3. Test of Goodness of fit analysis.

Statistics	Weibull	Gamma	Log-normal
Kolmogorov-Smirnov	0.051438	0.031019	0.041869
Anderson-Darling Criteria	359.574937	184.534433	177.549741
AIC	192915.1	191315.6	192463.5
BIC	192933.3	191333.7	192481.6

From table 3, it can be seen that two distributions fit the data meaning that the data. Using the Kolmogorov-Smirnov statistic it can be confirmed that the data followed the gamma and lognormal distributions better than the Weibull distribution since the two distribution statistics test values are lower than the critical value. Gamma fits the data best because from the Akaike's Information Criterion (191315.6) and Bayesian Information Criterion (191333.7) it is clearly evidenced that the distribution with smaller values is gamma distribution and therefore as per the decision rule it is considered the best of the three distributions for fitting the wind data. Also, it can be clearly observed that gamma is having the smallest Kolmogorov-Smirnov value (0.031019) which provides another evidence that the data follows gamma distribution best compared to the other distributions.

Table 4. Threshold values.

Threshold value	Distribution	AIC	BIC
0.1174	Weibull	190785.2	190803.3
	Gamma	190227.2	190245.3
	Log-normal	196888.9	196907
0.1180	Weibull	190777.2	190795.3
	Gamma	190228.9	190247
	Log-normal	197010.2	197028.3
0.1185	Weibull	190771.5	190789.6
	Gamma	190232	190250.1
	Log-normal	197133.1	197151.3
0.1190	Weibull	190766.7	190784.9
	Gamma	190237.3	190255.4
	Log-normal	197296.9	197315
0.1195	Weibull	190764.9	190783
	Gamma	190248.9	190267.1
	Log-normal	197561.8	197580
0.1199	Weibull	190773.6	190791.8
	Gamma	190279.	190298
	Log-normal	198192.7	198210.8

3.2. Minimum Distance Method for Three Parameters Distributions

First, there is need to investigate how the Akaike's Information Criterion and the Bayesian Information Criterion behaves under different value for the threshold parameter using 0.1174 as the baseline value from the data and then followed by several iterations. The threshold value with smaller AIC and BIC values was picked as the best. This will be investigated using table 4.

From table 4 generated using the distance method, it can be seen that gamma distribution has the smaller AIC and BIC value under all tested threshold values. This makes gamma the best distribution therefore we chose to use the threshold of 0.1174 for our analysis. This is because from the analysis of the original wind data it was found that gamma has a threshold value of 0.1174. Therefore, gamma is the best with the AIC value of 190227.2 and BIC value of 190245.3. Using the threshold value of 0.1174, the other two parameters for all the three distributions are given in table 5.

Table 5. Parameter estimation.

Distribution	Parameter	Estimate
Weibull	Shape	1.4115
	Scale	2.038586
Gamma	Shape	1.864567
	Scale	0.993702
Log-normal	Shape	0.4091521
	Scale	0.732395

To be sure that our data followed either of the three specific distributions, we performed a statistical test using the goodness of fit statistics. The goodness of fit statistics applied are summarized as shown in the table 6.

Table 6. Kolmogorov-Smirnov test.

Statistics	Weibull	Gamma	Log-normal
Kolmogorov-Smirnov	0.038007	0.028086	0.044478
Anderson-Darling	222.529802	120.255307	313.198226

From table 6, it can be confirmed using Kolmogorov-Smirnov that the data follows the gamma and Weibull distributions best than log-normal distributions. This is clearly supported by the fact that the test values for gamma and Weibull are smaller compared to the test value for lognormal distribution. Therefore, using statistical analysis it can be summarized that gamma three parameter distribution is the best among the three distributions for studying the Narok wind data.

Table 7. Best distribution using MDE.

Distribution	Criteria	Estimate
Gamma	AIC	190227.2
	BIC	190245.3
	Parameters	
	Threshold	0.1174
	Shape	1.864567
	Scale	0.993702

Using the minimum distance method, we therefore

conclude that gamma with three parameter distribution is the best with the following characteristics in table 7. This is because it has smaller AIC and BIC value compared to the gamma distribution with two parameters.

4. Conclusion

From the analysis, we conclude that using minimum distance estimation technique, gamma distribution with three parameters is the best distribution for fitting wind speed data

since it gives smaller AIC and BIC values. The scale parameter which shows how windy the site is (statistically the distribution of the wind speed) is estimated as 0.9937. The shape parameter which indicates how picked the site is (the most frequently expected wind speed) is estimated as 1.8646 and the threshold parameter which indicated the minimum expected wind speed value is estimated as 0.1174. Therefore, the gamma distribution with three parameters is the best distribution for investigating how windy or picked the region or site in Kenya is. The distribution is given as

$$f(y, z, q, t) = \frac{0.1174y^{1.8646-1}}{\Gamma(1.8646/0.1174)0.9937^{1.8646}} \exp \left[-\left(\frac{y}{0.9937} \right)^{0.1174} \right] \quad (19)$$

Where: Gamma function is treated as a continuous function depending on the wind speed data and y is the random variable (wind speed).

Author Contribution

O. O. K.; concept paper and proposal writing, O. O. K., O. E., A. A. D., and J. M.; study methods, O. O. K., O. E., and A. A. D.; data analysis (fitting of the probability distribution), O. E.; Validation.

Data Availability Statement

The data from the five sites is available at www.tahmo.org/climate.data.

Conflict of Interest

All the authors do not have any possible conflicts of interest.

Acknowledgements

The authors acknowledge School of Pure, applied and Health sciences at Maasai Mara University for offering ideological support during at all stages of this study.

References

- [1] Azami, Z., Khadijah, S., Mahir, A., and Sopian, K., (2009). Wind speed analysis in east coast of Malaysia. *European journal of scientific research*. Vol. 2.
- [2] Barasa, M., (2013). Wind regime analysis and reserve estimation in Kenya.
- [3] Celik, H., and Yilmaz, V., (2008). A statistical approach to estimate the wind speed distribution: the case study of Gelubolu region. Pp 122-132.
- [4] Galvao, F. A., and Wang, L., (2015). Efficient minimum distance estimator for quantile regression fixed effects panel data. *Journal of multivariate analysis*.
- [5] Gungor A. and Eskin, N., (2008). The characteristics that defines wind as an energy source.
- [6] Gupta, R., and Biswas, A., (2010). Wind data analysis of Silchar (Assam India) by Rayleigh and Weibull methods. *Journal of mechanical engineering research*. Vol. 2, pp 10-24.
- [7] Lawan, S. M., Abidin, W. A. W. Z., Chai, W. Y., Baharum, A., and Masri, T., (2015). Statistical modelling of long-term wind speed data. *American journal of computer science and information technology*.
- [8] Louzada, F., Ramos, P. L., and Gleici, S. C. P., (2016). Different estimation procedures for the parameters of the extended exponential geometric distribution for medical data. *Computational and mathematical methods in medicine*.
- [9] Lucen'o, A., (2006). Fitting the generalized pareto distribution to data using maximum goodness of fit estimators. *Computational statistics and data analysis*. Vol. 51. pp 904-917.
- [10] Mahyoub, H., (2006). Statistical anlysis of wind speed data and an assessment of wind energy potential in Taiz-Yemen. Vol. 2.
- [11] Maleki, F., and Deiri, E., (2007). Methods of estimation for three parameter reflected Weibull distribution.
- [12] Mert, I., and Karakus, C., (2015). A statistical analysis of wind speed using Burr, generalized gamma, and Weibull distribution in Antakya, Turkey. *Turkish journal of electrical engineering and computer science*.
- [13] Mumford, A. D., (1997). Robust parameter estimation for mixed Weibull (Seven parameters) including the method of maximum likelihood and the method of minimum distance. *Department of air force, Air force institute of technology*.
- [14] Oludhe, C., (1987). Statistical characteristics of wind power in Kenya. *University of Nairobi*.
- [15] Otieno, C. S., (2011). Analysis of wind speed based on Weibull model and data correlation for wind pattern description for a selected site in Juja, Kenya.
- [16] Otieno, F., Gaston, S., Kabende, E., Nkunda, F., and Ndeda, H., (2014). Wind power potential in Kigali and western provinces of Rwanda. *Asia pacific journal of energy and environment*. Vol. 1.
- [17] Rambachan, A., (2018). Maximum likelihood estimates and Minimum distance estimate.
- [18] Salma, O. B., and Abdelali A. E., (2018). Comparing maximum likelihood, least square and method of moments for Tas distribution. *Journal of Humanities and Applies science*.

- [19] Sanku, D., Menezes, A. F. B., and Mazucheli, J., (2019). Comparison of estimation methods for unit-Gamma distribution. *Journal of data science*. Vol. 17. pp 768-801.
- [20] Sukkiramathi, K., Sessaiah, C., and Indhumathy, D., (2014). A study of Weibull distribution to analyze the wind speed at Jogimatti in India. Vol. 01. pp 189-193.
- [21] Sultan, M. M. A., (2008). A data driven parameter estimation for the three parameter Weibull population from censored samples. *Mathematical and computational applications*. Vol. 13. Pp 129-136.
- [22] Ulgen, K., and Hepbasli, A., (2002). Determination of Weibull parameters for wind energy analysis of Izmir, Turkey.
- [23] Anderson, T. W., and Darling, D. A., (1952). Asymptotic theory of certain “goodness of fit” criteria based on stochastic processes. *The Annals of mathematical statistics*. Vol. 23, pp 193-212.