

Estimation of Missing Data and Maps of the Behaviour of Total Precipitation in Northern Colombia

Fernando Jove Wilches*¹, Rodrigo Hernández Avila¹ and Álvaro Rafael Caballero Guerrero²

¹ Department of Civil Engineering, Universidad de Sucre, Sincelejo, Sucre, Colombia.

² Department of Civil and Environmental Engineering, Universidad del Norte, Barranquilla, Colombia.

ORCIDs: 0000-0002-2080-4036 (Fernando), 0000-0003-3178-8075 (Rodrigo), 0000-0002-3567-9135 (Álvaro)

Abstract

The climatic conditions of a territory are a key element in the development of multiple types of projects, hence it is necessary to have relevant information that helps determine and simulate the future climate conditions that may influence such projects. At this point, the important thing is to have reliable and complete weather information records. Unfortunately, in Colombia, being able to find this type of information is not so simple, and in general the climatological data are incomplete for reasons such as failures or non-calibration in the measurement instruments, errors in the recording of measurements, meteorological extremes and difficulty in the access of the measurement areas. Hence, the determination of the missing data in the records of the climatological stations becomes a necessary question to complete the incomplete historical series. A common case is the missing data in rainfall stations, which are essential in hydrological simulation studies in order to obtain optimal hydraulic solutions. Therefore, it is almost mandatory to properly fill in missing data before performing any analysis on such data. This work presents the determination of missing data of total precipitation of 40 meteorological measurement stations located in the north of Colombia. The series analysed cover a period of 40 years of observations from 1980 to 2019, obtained by the Instituto de Hidrología, Meteorología y Estudios Ambientales (IDEAM). The estimation of missing data was carried out using Multiple Linear Regression (MLR), which allows estimating missing data based on nearby records that correlate with said data. The Beale – Little algorithm is applied which is based on the RLM method, using all the available information and leads to a simultaneous estimation of the missing values. The estimated precipitation dataset is useful for future hydrological analysis of studies that require meteorological forecasts in regions with losses of hydrological information, such as the meteorological stations located in the north of Colombia.

Keywords: Monthly rainfall, Missing data, Extreme rainfall, Multiple Linear Regression, Beale–Little algorithm.

I. INTRODUCTION

Hydrology is the science that deals with water, its occurrence, its circulation, distribution, its properties and its relationship with the environment and living beings [1]. Hydrology allows a rational and efficient use of the planet's water resources.

Through this, it is possible to plan and carry out studies that lead to the proper management of the water resources of a region.

Hydrology provides the technical and scientific elements that will allow you to adequately understand the water cycle [2]. These scientific methods generally require input data according to the models used and the specific conditions of the project. Rainfall is considered a main variable in hydrological studies, as it is the most important source for calculating the water balance and generating early warnings for risk of drought or flooding in a geographic region [3].

Rainfall analyses require continuous, homogeneous data that cover the maximum possible time interval. Unfortunately, the database available in the collection centres, most of the time, presents missing information due to the absence of reading, failure of the recording instrument, transcription error, among others, which limits the analyses and many times, it is constituted as a source of error [3]. Multiple Linear Regression (MLR) is a statistical technique that can be used in hydrology with the intention of determining the missing data on runoff characteristics, calculated in basins with hydrometric to sites or basins where such information is necessary and there are no gauges. [4]. In a complete way, the Beale-Little algorithm is used to simultaneously estimate the missing data in the records of the rainfall stations of the geographical area under study. With the use of this technique it is possible to obtain the data lost in the series of rainfall data, in order to carry out more reliable hydrological analyses and studies.

Once an area or region of the territory has been evaluated, the results of the analyses can be presented graphically. Being one of the most common and practical forms, show them in spatial distribution schemes where the hydrological variables analysed are represented. In the case of rainfall, the use of GIS tools for this purpose has become very popular. To carry out this, it is necessary to know the spatial variation in the study area, for which precipitation fields must be constructed, using interpolation methodologies, among which there is the Kriging method and the IDW method (Inverse Distance Weighted) [6]. Once the modelling has been carried out, the diagrams will be obtained showing the distribution and spatial interpolation of the rainfall in the study areas.

The study area presents different hydrological conditions throughout the entire territory, where there are geographically mountainous areas and others with swampy areas and river

discharge points, making their environmental conditions very varied. The information analysed includes precipitation data from 40 stations that have been active in the last 40 years and are located in the department of Sucre. The historical records of total rainfall were obtained from the data of the measurements of rainfall stations that rest in the Instituto de Hidrología, Meteorología y Estudios Ambientales (IDEAM).

The objective of this work is to present the RLM method through the use of the Beale-Little algorithm, for the determination of missing values of total rainfall in the stations located in the department of Sucre. From this, it was possible to fill in the missing rainfall for the forty stations in the study period, which was 40 years. In this way, it is possible to obtain reliable information that hydraulic researchers and designers require to develop their projects. Similarly, distribution maps of average and maximum rainfall of the study area were created. This allows the obtained results to be viewed in a clearer way and is of great use to professionals interested in studying the precipitation conditions of the area under study.

II. EXPERIMENTAL DESIGN, MATERIALS AND METHODS

II.I Study area description

The Department of Sucre is located in the Colombian Caribbean plain, north of the Central and Western mountain ranges, it has an area of 10,364 square kilometres. It limits to the north and east with the Department of Bolívar, to the south with the Departments of Antioquia and Córdoba, to the west with the Department of Córdoba and to the northeast with the Caribbean Sea. A little more than a third of its territory forms the flood-prone depression of the Bajo Magdalena, Cauca and San Jorge rivers, characterized by numerous swamps, especially along the San Jorge river. To the northwest, on the other hand, a strip of hills can be seen, corresponding to the San Jacinto or Montes de María mountains. Between these hills and the San Jorge depression lie more or less flat savannas. In the Department of Sucre the climate is warm, dry towards the sea and humid towards the depression. The various reliefs in the Department of Sucre have a warm thermal floor with temperatures that range between 25.5 °C and 28.7 °C average per year. The precipitation regime is determined by the geographical location, and by the influence of some factors, such as atmospheric circulation, relief, the interaction between land and sea and the influence of jungle or wooded areas. In Sucre, the annual average rainfall varies between 1,000 mm for the less humid areas in the north, up to 2,800 mm in the rainiest areas of the south [7].

II.II Material and methods

The determination of the missing rainfall data was carried out by applying the Beale-Little algorithm, which is a Multivariate Analysis (AM) technique that allows to simultaneously estimate the missing data in records of hydrometric or pluviometric stations of a geographical area, which show a significant correlation, but do not have persistence. The technique is appropriate to estimate missing data that were lost

randomly, that is, when there was an absence of the operator or a temporary suspension, maintenance or poor calibration of the equipment, or, due to improvements in the installation [8] [9].

The first step in the methodology is to collect the available information and organize it in a matrix arrangement where the rows represent the recording periods and the columns represent the meteorological stations, adding an asterisk to the missing data or leaving the spaces in blanks. Next, the initial substitution of the missing data is made, for this, the arithmetic average of each column is calculated and this value is replaced in the spaces with missing data. From these, the MLR for each record is calculated and the missing data is estimated. After applying the RLM, the differences between the old data and the new data are calculated, which would constitute the first cycle of the Beale-Little algorithm. This process continues for the number of cycles that is necessary until the difference between the new data and the previous one takes a sufficiently small value or is almost negligible; so that if this degree of precision is not reached, the new estimates are replaced by the old ones and new estimates are made using the same procedure.

The mathematical formulation of the RLM [10] [11] is shown below. It is assumed that there are n observations of Y , X_1 , X_2, \dots, X_p , considering that there are p independent variables or regressors of the Regression, whose model is the following:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon \quad (1)$$

In matrix notation, we have:

$$Y = X \cdot \beta + \varepsilon \quad (2)$$

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1p} \\ 1 & X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \dots & \dots & \dots & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

The matrix X contains the X_{ij} observations of the meteorological data, showing the i -th periods of the observations and the j -th independent variables. The method seeks that the sum of the squared errors (ε_i) is minimized according to the expression:

$$\sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \beta_2 X_{i2} - \dots - \beta_p X_{ip})^2 \quad (3)$$

When differentiating the right hand side of equation 3, based on $\beta_0, \beta_1, \dots, \beta_p$ separately and equaling to zero, this produces p equations with p unknown parameters, and is called normal equations. This, written in matrix notation, looks like this:

$$X^T \cdot X \cdot \beta = X^T \cdot Y \quad (4)$$

The solution of equation 4 for the β_i can be expressed from Equation 5:

$$\beta = (X^T \cdot X)^{-1} \cdot X^T \cdot Y \quad (5)$$

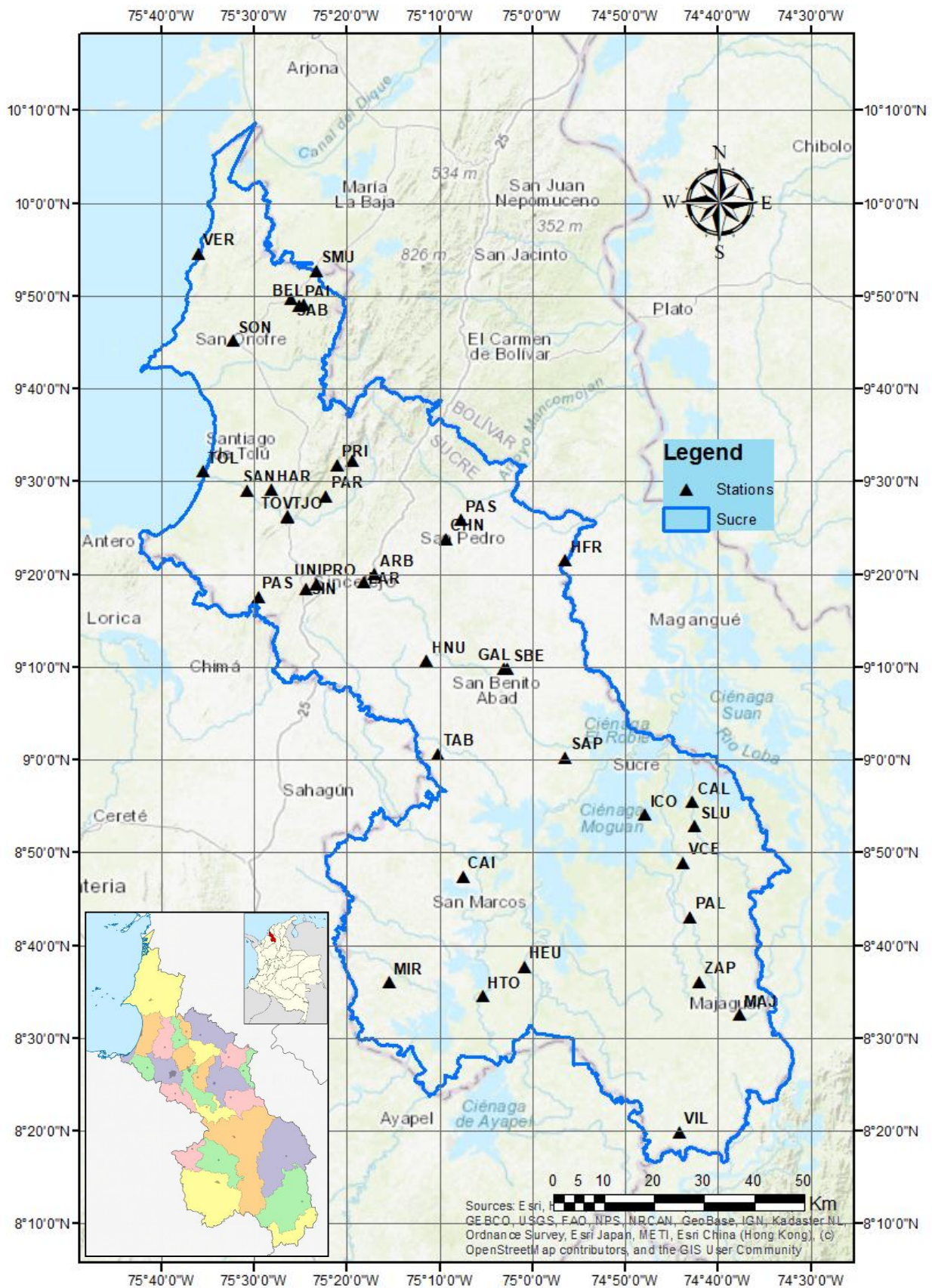


Fig 1. Location of meteorological stations belonging to the department of Sucre

In equation 5, X is the matrix with the observations and variables, Y is the vector of dependent observations, X^T is the transpose matrix of X , and $(X^T \cdot X)^{-1}$ is the matrix inverse of the matrix constructed through $X^T \cdot X$.

Additionally, Figures were created showing the average and maximum rainfall of the stations, in order to generate an extrapolation of the data and observe the behavior of rainfall based on the historical series with missing data filled in, being able to observe the behavior of average rainfall over the 40-year Study period.

To carry out this process, the GIS computing tool, ArcMap, was used, for which the geographic location data of the stations were imported into the program, as well as the values related to the average and maximum precipitation. Additionally, layers were imported to limit the study area (the department of Sucre) and other layers to give relief to neighboring places and general information.

The Kriging method was used to perform the spatial interpolation of the precipitation variable. This is a geostatistical method, which relates mathematics to earth sciences, therefore being ideal for estimating the rainfall data

obtained from the measurement of the rain gauge. The Kriging method works with statistical models that involve autocorrelation [6].

Autocorrelation allows spatial interpolation to occur, which generally occurs when estimating a regionalized value at unsampled points, based on the weight of the regionalized values that were observed. Equation 6 shows the mathematical model that governs spatial interpolation [6].

$$Zg = \sum_{i=1}^{ns} \lambda_i Zsi \quad (6)$$

In equation 6, Zg is the interpolated value in the requested areas or points, Zsi is the value observed at point i , ns is the total number of points observed and λ_i is the weight that contributes to the interpolation.

Within ArcMap this method (Kriging) was selected, because it has some advantages over the traditional interpolation methods built into GIS. If the assumptions of the theory hold, the method is less arbitrary compared to others. A very important difference to emphasize is that the weights used in the method are determined by the variogram and the configuration of the observed data [12].

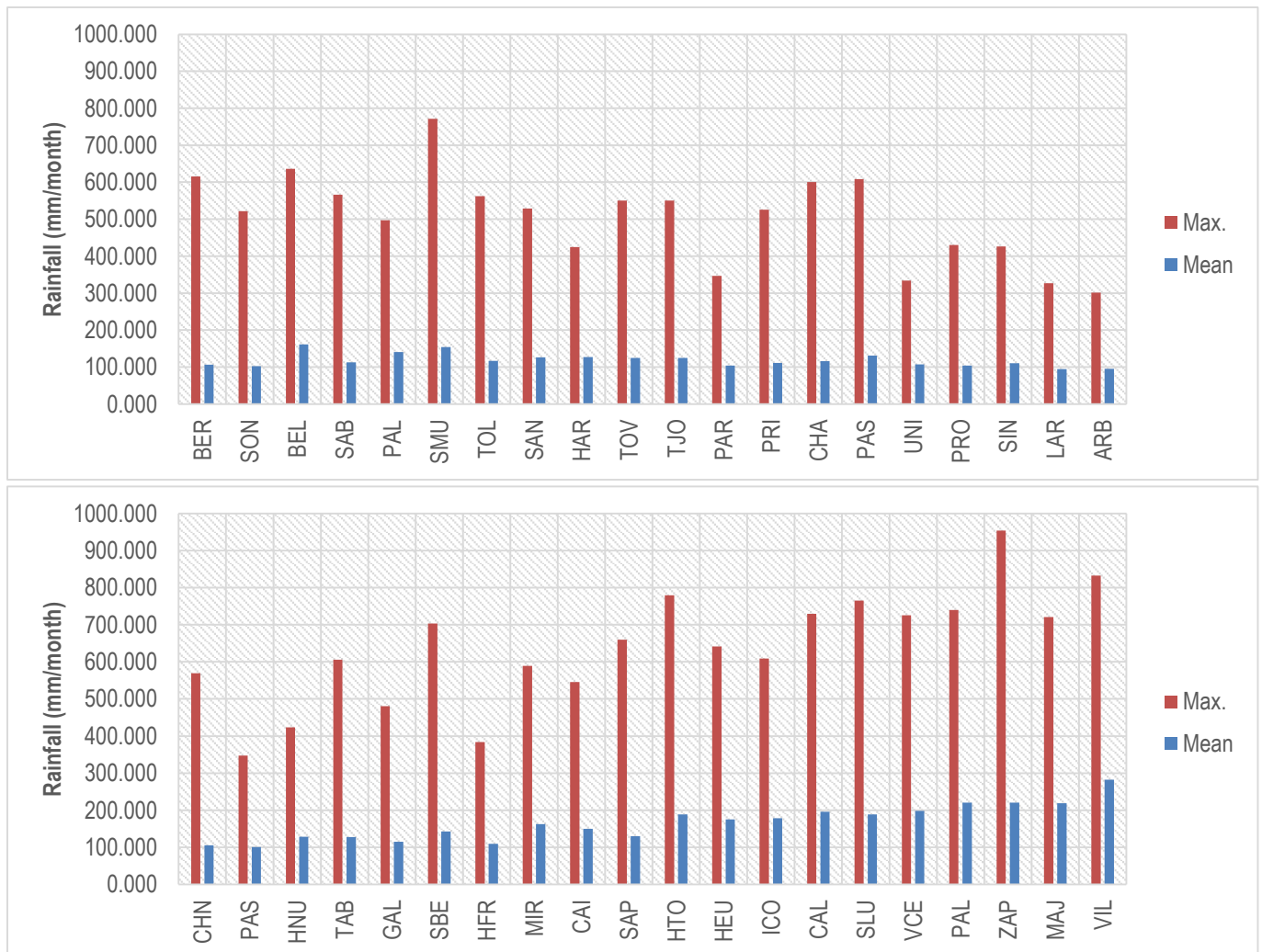


Fig. 2. Mean and maximum values of monthly rainfall.

Table 1. General data of the Meteorological Stations

ID	NUMBER	STATION NAME	CODE	ALTITUDE	LATITUDE	LONGITUDE	MUNICIPALITY
1	13090030	BERRUGAS	BER	1	9.91027778	-75.6005556	San Onofre
2	13090020	SAN ONOFRE	SON	55	9.75416667	-75.5380556	San Onofre
3	13090080	BELEN HACIENDA	BEL	60	9.81944444	-75.4102778	San Onofre
4	13090010	SABANETICA	SAB	1	9.82888889	-75.4327778	San Onofre
5	13090110	PALO ALTO	PAL	50	9.81750000	-75.4183333	San Onofre
6	29030380	SABANAS DE MUCACAL	SMU	10	9.87980556	-75.3886944	San Onofre
7	13090070	TOLU	TOL	2	9.51916667	-75.5911111	Tolú
8	13090090	SANTA ANGELA HACIENDA	SAN	20	9.48527778	-75.5125000	Tolú
9	13090100	ARGENTINA LA HACIENDA	HAR	20	9.48583333	-75.4686111	Tolú Viejo
10	13090050	TOLUVIEJO	TOV	60	9.43800000	-75.4400000	Tolú Viejo
11	13097010	TOLUVIEJO	TJO	59	9.43700000	-75.4390000	Tolú Viejo
12	13090060	PARAISO EL	PAR	100	9.47388889	-75.3708333	Colosó
13	13095020	PRIMATES	PRI	200	9.53013889	-75.3513611	Colosó
14	13090040	CHALAN	CHA	100	9.53861111	-75.3225000	Chalán
15	13090230	PASTORA LA	PAS	58	9.29333333	-75.4919444	Sincelejo
16	25025270	UNISUCRE	UNI	160	9.31638889	-75.3875000	Sampués
17	25025350	PUERTA ROJA	PRO	160	9.31638889	-75.3875000	Sincelejo
18	25020130	SINCELEJO	SIN	200	9.30816667	-75.4068333	Sincelejo
19	25020400	LIBRA ARRIBA	LAR	180	9.32055556	-75.3036111	Betulia
20	25025080	APTO RAFAEL BARVO	ARB	166	9.33388889	-75.2830556	Corozal
21	25020390	CHARCON	CHN	150	9.39722222	-75.1558056	Los Palmitos
22	25020190	PUERTO ASIS	PAS	200	9.43194444	-75.1291667	Los Palmitos
23	25020990	HATO NUEVO	HNU	80	9.17833333	-75.1902778	Corozal
24	25020750	TABLITAS LAS	TAB	60	9.01138889	-75.1688889	San Benito Abad
25	25021000	GALERAS	GAL	70	9.16500000	-75.0508333	Galeras
26	25025380	SAN BENITO	SBE	20	9.16388889	-75.0447222	San Benito Abad
27	25020860	FRONTERA LA HACIENDA	HFR	100	9.35944444	-74.9408333	Buenavista
28	25021660	MIRASOL	MIR	30	8.60083333	-75.2572222	San Marcos
29	25020980	CAIMITO	CAI	20	8.79083333	-75.1244444	Caimito
30	25020760	SANTIAGO APOSTOL	SAP	25	9.00472222	-74.9405556	San Benito Abad
31	25021470	TORNO EL HACIENDA	HTO	60	8.57638889	-75.0897222	San Marcos
32	25020740	EUREKA HACIENDA	HEU	20	8.62916667	-75.0147222	San Marcos
33	25021560	ISLA DEL COCO	ICO	20	8.90277778	-74.7986111	Sucre (Sucre)
34	25021360	CAMPO ALEGRE	CAL	20	8.92472222	-74.7119444	Sucre (Sucre)
35	25021370	SAN LUIS	SLU	20	8.88166667	-74.7080556	Sucre (Sucre)
36	25020500	VILLA CECILIA HACIENDA	VCE	50	8.81611111	-74.7294444	Sucre (Sucre)
37	25020790	PALMARITO	PAL	50	8.71888889	-74.7177778	Majagual
38	25020820	ZAPATA	ZAP	50	8.60277778	-74.6997222	Majagual
39	25025240	MAJAGUAL	MAJ	26	8.54269444	-74.6273333	Majagual
40	25020940	VILLANUEVA	VIL	45	8.33277778	-74.7355556	Guaranda

Table 2. Statistical Description of the Data of the Meteorological Stations

CODE	Observations	Minimum (mm/month)	Maximum (mm/month)	Mean (mm/month)	Standard deviation (mm/month)
BER	480	0.161	615.500	106.155	103.801
SON	480	0.554	521.600	102.349	83.100
BEL	480	0.100	636.208	161.092	113.493
SAB	480	0.080	566.000	113.084	98.949
PAL	480	0.233	497.167	140.471	97.700
SMU	480	0.304	771.125	154.512	113.475
TOL	480	0.277	562.500	116.658	101.510
SAN	480	0.708	528.667	126.011	104.181
HAR	480	0.312	425.017	127.547	95.134
TOV	480	0.100	550.000	124.710	98.335
TJO	480	0.055	550.000	125.031	98.848
PAR	480	0.400	346.742	103.775	75.258
PRI	480	0.736	525.858	111.484	89.300
CHA	480	0.200	600.375	116.309	94.523
PAS	480	0.040	608.417	131.407	103.449
UNI	480	0.273	334.194	107.471	75.653
PRO	480	1.000	430.000	104.139	72.846
SIN	480	1.220	426.000	110.153	77.261
LAR	480	0.896	326.690	94.822	71.803
ARB	480	0.088	301.700	95.164	68.821
CHN	480	0.468	569.500	105.440	80.147
PAS	480	0.175	347.375	100.986	70.993
HNU	480	1.000	423.208	128.674	96.774
TAB	480	0.500	605.500	127.729	98.497
GAL	480	0.331	480.667	115.122	92.116
SBE	480	0.146	703.158	143.126	113.750
HFR	480	0.583	383.583	109.328	79.578
MIR	480	0.723	589.000	162.386	126.330
CAI	480	0.430	545.375	150.230	117.609
SAP	480	0.200	660.000	130.086	109.792
HTO	480	0.391	779.792	188.760	154.414
HEU	480	1.000	641.792	175.472	134.026
ICO	480	1.458	608.875	178.659	137.756
CAL	480	0.886	729.375	195.968	155.168
SLU	480	0.898	765.417	188.797	150.821
VCE	480	0.106	725.292	198.295	155.765
PAL	480	0.036	740.000	220.580	175.848
ZAP	480	0.997	954.042	220.417	169.662
MAJ	480	0.700	720.917	218.938	163.761
VIL	480	1.417	832.167	282.592	199.068

III. RESULTS

This section contains Figures and Tables of the precipitation data that were analyzed from the information obtained from 40 meteorological stations located within the territory of the Department of Sucre (Colombia).

Fig. 1 shows the map of the location of the department of Sucre and the location of the meteorological stations that were selected for this work. The names, codes, the geographic coordinates of the location, the elevations and the municipalities where the IDEAM meteorological stations are located were tabulated in Table 1. The information on the

statistical description of the data of the selected meteorological stations is shown in Table 2, where you can see the nomenclature given to each station, the number of observations considered during the study period (40 years), as well as the maximum, average values and the corresponding standard deviation, related to the Total monthly rainfall for each station.

Fig. 2 shows the historical average of the mean and maximum annual rainfall in the 40-year period for the department of Sucre. The stations are organized in a North-South direction, thus allowing us to look at their behavior.

In Fig. 3 and 4, a graphical representation has been made through the ArcMap program of the average and maximum precipitations of the area under study, in this way it can be easily appreciated which are the parts of the department where it has rained the most in the last 40 years and which have been the places with less precipitation, which gives an idea of the spatial and temporal distribution of the rains during the study period.

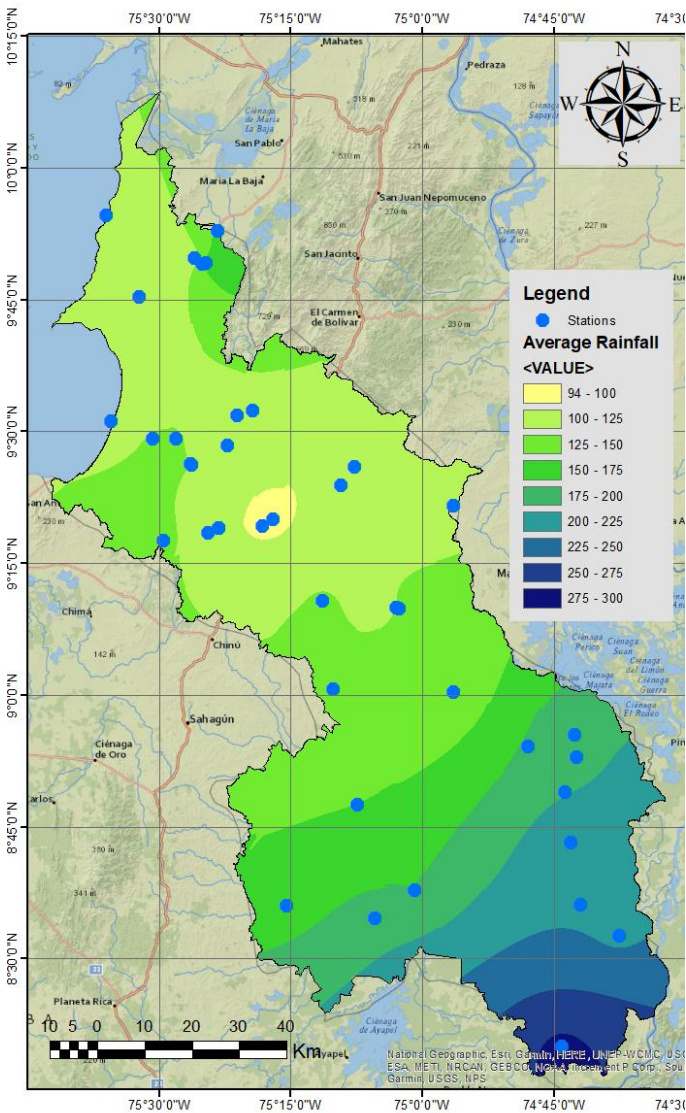


Fig. 3. Comportamiento de la precipitación promedio mensual en el departamento de Sucre

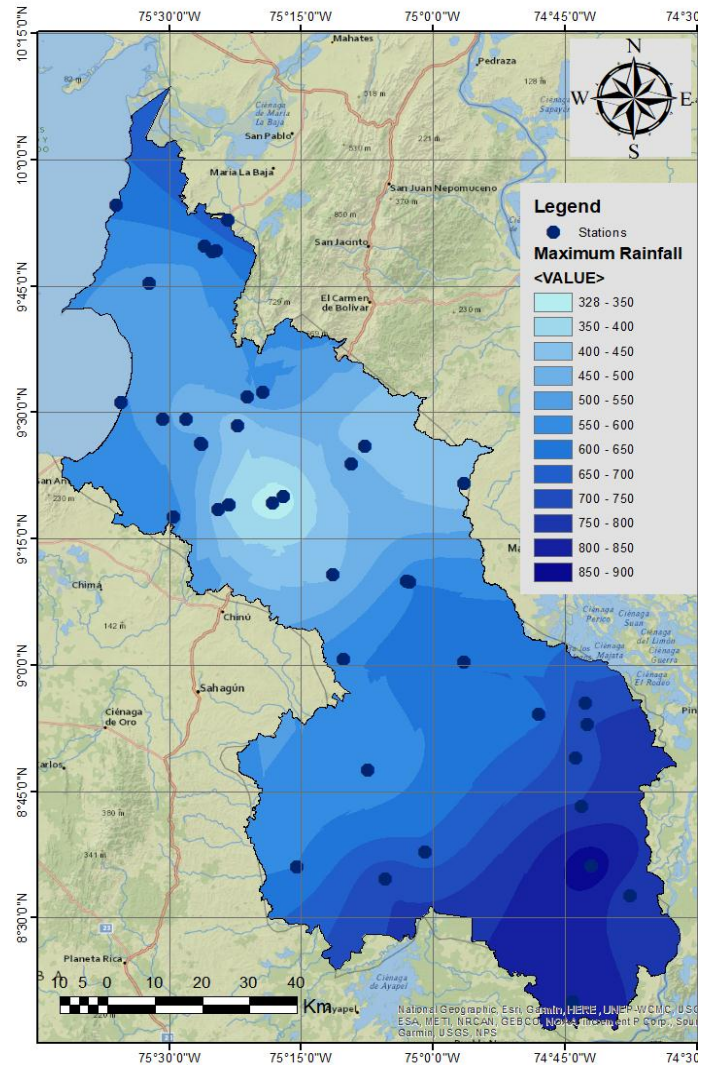


Fig. 4. Comportamiento de la precipitación máxima mensual en el departamento de Sucre

IV. CONCLUSIONS

From the work carried out, the missing monthly precipitation data could be obtained for the 40 selected meteorological stations. This information is of particular interest to be taken as an essential input for the hydrological characterization of the engineering projects that are intended to be developed within the department of Sucre.

Additionally, maps were obtained that show the behavior of monthly average and maximum monthly rainfall in the department of Sucre, information that allows establishing the rainfall characteristics of a project, according to its geographic location.

Acknowledgments

The authors thank the Instituto de Hidrología, Meteorología y Estudios Ambientales (IDEAM) for providing the database containing the monthly and annual rainfall for historical periods of all the stations in the department of Sucre

(Colombia). Authors also thank the engineer Edilberto Elias Contreras Sierra for the contributions made in the development of the manuscript.

REFERENCES

- [1] Linsley R, Kohler M, Paulhus J. Hidrología para ingenieros, 2 ed. New York, McGraw-Hill, 1985.
- [2] Legarda Burbano L, Viveros Zarama M. La importancia de la hidrología en el manejo de cuencas hidrográficas. *Revista De Ciencias Agrícolas*.1996;14(1y2). Disponible en <https://revistas.udenar.edu.co/index.php/rfacia/article/view/1163?articlesBySameAuthorPage=3#articlesBySameAuthor>
- [3] Herrera C, Campos J, Carrillo J. Estimación de datos faltantes de precipitación por el método de regresión lineal: Caso de estudio Cuenca Guadalupe, Baja California, México. *Investigación y Ciencia*. 2017; 25 (71): 34-44. Disponible en <https://www.redalyc.org/articulo.oa?id=67452917005>
- [4] Campos-Aranda D. Una aplicación hidrológica de la regresión lineal múltiple ponderada. *Tecnología y ciencias del agua*. 2016; 7 (4): 161-173. Disponible en http://www.scielo.org.mx/scielo.php?script=sci_arttext&pid=S2007-24222016000400161&lng=es&nrm=iso
- [5] Vargas A, Santos A, Cárdenas E, Obregón N. Análisis de la distribución e interpolación espacial de las lluvias en Bogotá, Colombia. *Dyna*. 2011; 78 (167): 151-159. Disponible en: <https://www.redalyc.org/articulo.oa?id=496/49622358017>
- [6] Ly S, Charles C, Degré A. Different methods for spatial interpolation of rainfall data for operational hydrology and hydrological modeling at watershed scale: a review, *BASE*. 2013; 17 (2).
- [7] Banco de la Republica de Colombia, Documentos de trabajo sobre la economía regional. https://www.banrep.gov.co/docum/Lectura_finanzas/pdf/DTSER-63-VE.pdf, 2005 (accessed 9 September 2020)
- [8] Beale EML, Little RJA. Missing values in multivariate analysis, *Journal of Royal Statistical Society B*. 1975; 37: 129-145. DOI; <https://doi.org/10.1111/j.2517-6161.1975.tb01037.x>
- [9] Clarke RT. Statistical modelling in hydrology. Capítulo 7, *Multivariate models*, pp. 254-302, Chichester, Inglaterra, John Wiley & Sons, 1994. 412.
- [10] Ryan TP. Linear Regression, capítulo 14, pp. 14.1-14.43, en: *Handbook of Statistical Methods for Engineers, Scientists*, Harrison, M., Wadsworth, editor, 2a. ed., Nueva York, McGraw-Hill Book Co., 1998.
- [11] Campos-Aranda D, Estimación simultánea de datos hidrológicos anuales faltantes en múltiples sitios, *Ingeniería, Investigación y Tecnología*. 2015; 16(2): 295-306. <https://doi.org/10.1016/j.riit.2015.03.013>
- [12] Oliver MA, Kriging: A Method of Interpolation for Geographical Information Systems, *International Journal of Geographic Information Systems*. 1990; 4: 313–332.