# Probabilistic Bayesian Analysis and Inference using Geophysical well-logging
# Case Four: Transit Time Data with Sonic Log and Gamma Log

## M. Sc. Zenteno Jiménez José Roberto

*Geophysical Engineering, National Polytechnic Institute, Mexico City, ESIA-Ticóman Unit Mayor Gustavo A. Madero*
jzenteno@ipn.mx

**Abstract:** The methodology was used to obtain new normal probability and extreme value distribution functions through Bayesian inference and stochastic mixture of Gaussians. The proposed methodology is oriented to data with Gaussian behavior and consists of fitting the normal distribution function to the time series data of the porosity data, transit time curves and the curve of the gamma geophysical record, to find the probability of the Minimum or estimated values to find areas of interest, then we use Bayesian inference for normal data, in this case they are looking for behavior and trend with new Gaussian and extreme value functions. To validate the model we use the following statistical estimators, measurement of the root of the squared error, squared error, coefficient of determination and prediction approximation.

**Keywords:** Bayesian Inference, Gaussian Mixing, Extreme Variable Distribution Functions.

## Introductión

In the normal distribution, equation (1) one can calculate the probability that various values occur within certain ranges or intervals. However, the exact probability of a particular value within a continuous distribution, such as the normal distribution, is zero. This property distinguishes continuous variables, which are measured, from discrete variables, which are counted. As an example, time (in seconds) is measured and not counted.

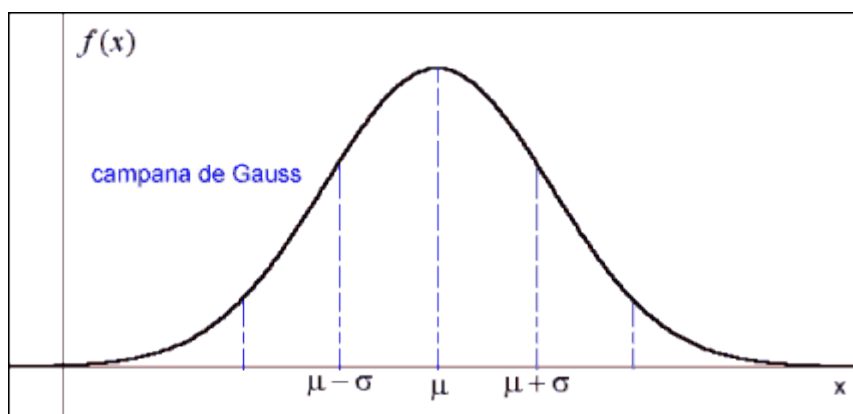$$\varphi(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{\left(-\frac{(x-\mu)^2}{2a^2}\right)} \tag{1}$$



Figure 1. Gaussian bell or Gaussian density function (Source: Internet)

Now as we obtain the estimation parameters of a probability distribution function, we will use the Maximum Likelihood technique for the parameters to be estimated for the fit. The maximum likelihood method is a procedure to obtain a point estimator of a random variable.

Let $(X1,…, Xn)$ be a random sample with a distribution function $f(x \mid \theta)$.

We define the likelihood function as:

$$L(\theta|X_1, X_2, \dots X_n) = \prod_{i=1}^{n} f(X_i|\theta) \tag{2}$$

The estimator of θ in the maximum likelihood method is the value that maximizes the likelihood function. This value is called the maximum likelihood estimator EMV (θ).

Now:

$$L(\theta|X_1, X_2, \dots X_n) = \ln\big(L(\theta|X_1, X_2, \dots X_n)\big) = \sum_{i=1}^{n} f(X_i|\theta) \tag{3}$$

Therefore, the maximum likelihood estimator is defined as:

$$EMV(\theta) = \max_{\theta \in \theta} L(\theta|X_1, X_2, \dots X_n) \tag{4}$$

Through the previous description, the parameters of a Normal Distribution are obtained, which are by doing the following:

By the Maximum Likelihood Method, the likelihood function will be:

$$L(\mu, \sigma) = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^N e^{\left(-\frac{\sum(x-\mu)^2}{2\sigma^2}\right)} \tag{5}$$

Applying logarithms and deriving with respect to the parameters to be estimated, we have a system of equations as follows:

$$\frac{\partial \log(L(\mu, \sigma))}{\partial \mu} = \frac{\sum x}{\sigma^2} - \frac{N_\mu}{\sigma^2} = 0 \tag{6}$$

$$\frac{\partial \log(L(\mu, \sigma))}{\partial \sigma} = -\frac{N}{\sigma} + \frac{\sum(x-\mu)^2}{\sigma^3} = 0$$

With the solution:

$$\mu = xy\sigma^2 = \frac{1}{N}\sum_{i=1}^{N}(x-\mu)^2 \tag{7}$$

With the first adjusted Normal we get

$$Normal(\mu, \sigma^2) \tag{8}$$

Now what we are looking for are extreme values that we want to know, the probability of occurrence so we use Bayesian Inference to find this probability with a new distribution function that will be part of the new functions of normal distribution and function of extreme variable or gev, our new unknown will be the mean, we can also proceed by observing the behavior for this case of a very important petrophysical variable: Porosity

Porosity is an important measure since what we are looking for is the saturation of the fluid found in the formation or rock, there are various methods to find analytical methods such as formulas, through direct measurement in the oil well or with a direct petrophysical analysis of the core or a sample of the formation.

### Bayesian inference
Bayesian inference is the process of analyzing statistical models with the incorporation of prior knowledge about the model or the model parameters. The root of such an inference is Bayes' theorem:

$$P(Parameters|Data) \tag{9}$$
$$= \frac{P(Data|Parameters) * P(Parameters)}{P(Data)}$$
$$\approx FVerosimilitud * PDF\ Priori$$

In this case we have the observations in the normal distribution form

$$X|\theta \sim N(\theta, \sigma^2) \tag{10}$$

Where the sigma is previously known and the PDF a Priori is

$$\theta \sim N(\mu, \tau^2) \tag{11}$$

Here mu and tao are also known, we are looking for n samples of the observed data, in the case of Ozone the maximum values or above 150 ppb, the Case of the particulate PM10 above 120 microgr / m3, the Case of PM2. 5 above 65 microgr / m3 and in the Case of Maximum Temperatures is the entire Data sample and thus We Obtain the New Normal Distribution Function with the new searched parameter, now for this new case the ones we are going to try to find is the zone with the lowest value of porosities with this stochastic approximation, as we proceed previously in previous works with Bayesian inference:

$$\theta|X \sim NB\left(\frac{\tau^2}{\frac{\sigma^2}{n} + \tau^2} * X + \frac{\frac{\sigma^2}{n}}{\frac{\sigma^2}{n} + \tau^2} * \mu, \frac{\frac{\sigma^2}{n} * \tau^2}{\frac{\sigma^2}{n} + \tau^2}\right) \tag{12}$$

Now these data contain noise, there are values very close to zero, from the adjustment process then we proceed to adjust the petrophysical variable such as porosity and if it has an adjustment as such to the Gaussian is more favorable, we proceed to count the number of low values of the which we know could be, depending on the rock and the behavior of the Geophysical Record in this case which will be our mean $\tau$ and our $\sigma$ as the standard deviation found with the adjustment, we apply a function to find the minimum value and we introduce a random Gaussian with that number of low values found the minimum as a vector, the $\sigma$ the mean the number of values below the estimate from the read record, the mean $\mu$, will be the mean of the complete porosity time series of the data and $\tau$ as the same $\sigma$

Subsequently we apply a random noise with a Uniform Distribution with the length of the terms of the time series of the data, we now apply an adjustment with the Extreme Value Distribution Function (GEV) to find the parameters that fit even better these Data with the random noise and thus find the distribution functions that will give us the possible zones.

$$GEV(\mu, \sigma, k) \tag{13}$$

A GEV is fitted with this uniform random distribution

$$GEVa(Xa, \mu, \sigma, k1) \tag{14}$$

And later a GEVA is generated (with the random parameters of GEVa)

$$GEVA(\mu a, \sigma a, k1) \tag{15}$$

The extreme value functions or GEVs with the new parameters of the new extreme distribution functions see [14]:

$$GEV\left(\sum_{i=1}^{n} \frac{\mu_i}{n}, \frac{1}{n-1}\sum_{i=1}^{n-1} \sigma_i, \sum_{i=0}^{n} \frac{k_i}{n}\right) \tag{16}$$

We have the following probability distribution functions, the Bayesian Normal, GEV of the data and the Random GEV, therefore they are 3 functions that we have so the first two sums of two probability distribution functions.

**Table 2. GEV and Normal probability distribution functions**

| GEV | Normal |
|---|---|
| $GEV(\mu, \sigma, k)$ | $Normal(\mu, \sigma^2)$ |
| $GEV(Mupostmean\ \mu, SigmaSD\ \sigma, k2)$ | $Normal1(E(GEV1), \sqrt{Var(GEV1)})$ |
| $GEV(Mupostmean2\ \mu, SigmaSD2\ \sigma, k2)$ | $Normal2(E(GEV2), \sqrt{Var(GEV2)})$ |

### Adjustment Indicators

Indicators of deviation of a group of data in relation to a model can be used to assess the goodness of fit between the two. Among the most common indicators are the following. Those that were used to determine the distribution that best fit the data. They are the mean square error (RMSE), mean square error (MSE), prediction precision (PA) and coefficient of determination (R2) Table 4 gives the equations for the adjustment indicators that have been used by Lu (2003) and Junninen et al. (2002).

**Table 3.**

| Indicator | Equation |
|---|---|
| **Root Mean Square Error (Raíz Cuadrada del Error)** | $RMSE = \sqrt{\left(\dfrac{1}{N-1}\right) \sum_{i=1}^{N} (Pi - Oi)^2}$ |
| **Mean Square Error (Error Cuadrado Principal)** | $MSE = \left(\dfrac{1}{N}\right) \sum_{i=1}^{N} (Pi - Oi)^2$ |
| **Coeficiente of Determination (Coeficiente de Determinación)** | $R^2 = \left(\dfrac{\sum_{i=1}^{N}(Pi - P)(Oi - O)}{N S_p S_o}\right)^2$ |
| **Prediction Accuracy (Precisión de Predicción)** | $AP = \dfrac{\sum_{i=1}^{N}(Pi - O)^2}{\sum_{i=1}^{N}(Oi - O)^2}$ |

Notation: N = Number of Observations, Pi = Predictive Values, Oi = Observed Values, P = Average of Predicted Values, O = Average of Observed Values, Sp = Standard Deviation of Predicted Values, So = Standard Deviation of Values Observed.

**Table 4. Fit Indicators for Each Fitted Gaussian**

| Gaussian Porosity | RMSE | MSE | $R^2$ | AP | K S Test |
|---|---|---|---|---|---|
| Well 4 | 0.1741 | 0.0303 | 0.9654 | 0.5067 | 5.2281e-04 |
| Well 2 | 0.3574 | 0.7947 | 0.8916 | 1.3250 | 0.0079 |
| Well 3 | 0.2244 | 0.0503 | 0.9467 | 0.8536 | 0.0 |
| Gaussian GR | RMSE | MSE | $R^2$ | AP | K S Test |
| Well 3G | 0.2442 | 0.0596 | 0.9518 | 0.4397 | 0.0 |

### Stochastic Method of Gaussian Mixtures

The clustering model most closely related to statistics is the distribution-based model. The groups can then easily be defined as the objects that most likely belong to the same distribution. A convenient property of this approximation is that it is very similar to the way artificial data sets are generated: by random sampling of objects from a distribution.

One of the most prominent methods is known as the Gaussian mixture model (used in the expectation-maximization algorithm). Here, the data set is normally modeled with a fixed number (to avoid overfitting) of Gaussian distributions that is randomly initialized, and whose parameters are iteratively optimized to better classify the data set. This will converge to a local optimum, multiple runs can produce different results. To obtain a good grouping, the data are often assigned to the Gaussian distribution with the highest probability of belonging to such a grouping.

Gaussian mixture models are a probabilistic model for representing normally distributed subpopulations within a general population. Mixture models in general do not require knowing which subpopulation a data point belongs to, allowing the model to automatically learn the subpopulations, using Expectation-Maximization (EM).

A Gaussian mixture model means that each data point is put (randomly) from one of the data classes C, with probability p_i of being drawn from class i, and each class is distributed as Gaussian with mean standard deviation μ_i and σ_i. Given a set of data extracted from said distribution, we seek to estimate these unknown parameters.

The algorithm used here for estimation is EM (Expectation Maximization). In short, if we knew the class of each of the N input data points, we could separate them, and use Maximum Probability to estimate the parameters of each class. This is the step that makes (soft) selections of (unknown) classes for each of the data points based on the previous round of parameter estimates for each class.

$$
\begin{aligned}
\phi_j &:= \frac{1}{m}\sum_{i=1}^{m} w_j^{(i)}, \\
\mu_j &:= \frac{\sum_{i=1}^{m} w_j^{(i)} x^{(i)}}{\sum_{i=1}^{m} w_j^{(i)}}, \\
\Sigma_j &:= \frac{\sum_{i=1}^{m} w_j^{(i)}\big(x^{(i)}-\mu_j\big)\big(x^{(i)}-\mu_j\big)^{T}}{\sum_{i=1}^{m} w_j^{(i)}}
\end{aligned}
$$

Figure 2. Basic equations of the EM Algorithms (Source:http://mccormickml.com/2014 )

## Results

**Well 4 - Porosity and Transit Time**



**Inference**

*International Journal of Latest Research in Engineering and Technology (IJLRET)*
*ISSN: 2454-5031*
*www.ijlret.com || Volume 06 - Issue 09 || September 2020 || PP. 13-31*

**Well 2 Apertura (Porosity case of Gaussian skewed to the right)**

**Inference**

*International Journal of Latest Research in Engineering and Technology (IJLRET)*
*ISSN: 2454-5031*
*www.ijlret.com || Volume 06 - Issue 09 || September 2020 || PP. 13-31*

**Well 3– Porosity and Transit Time**



**Figure 3. Well Log 3**
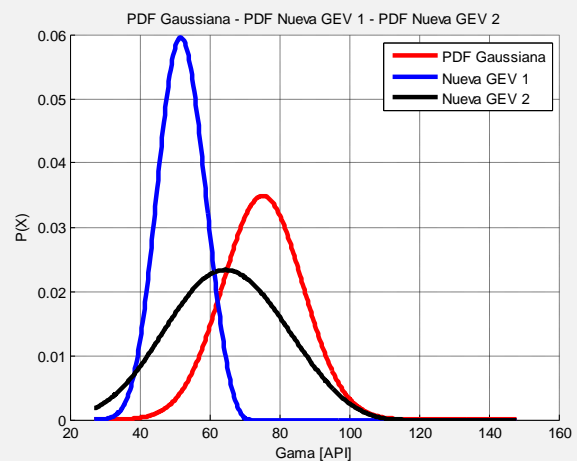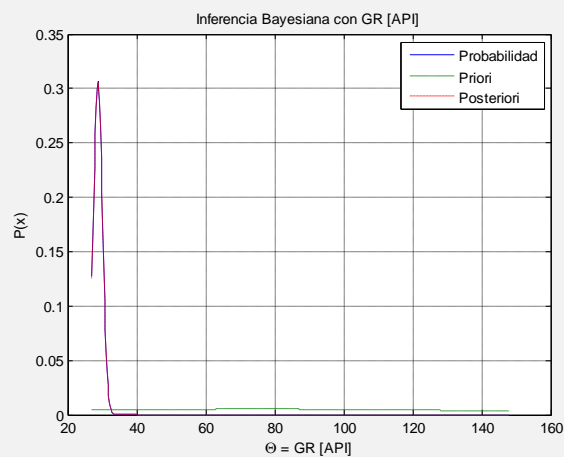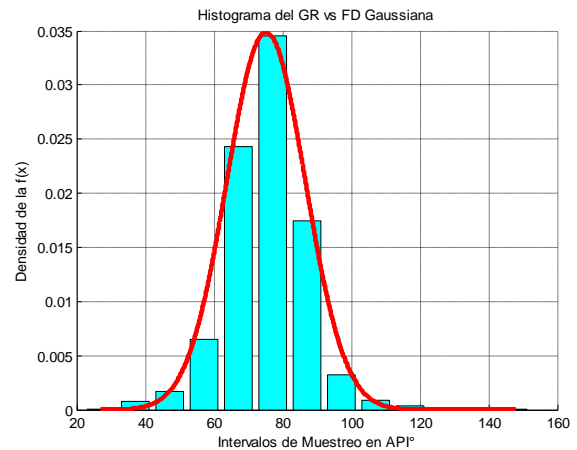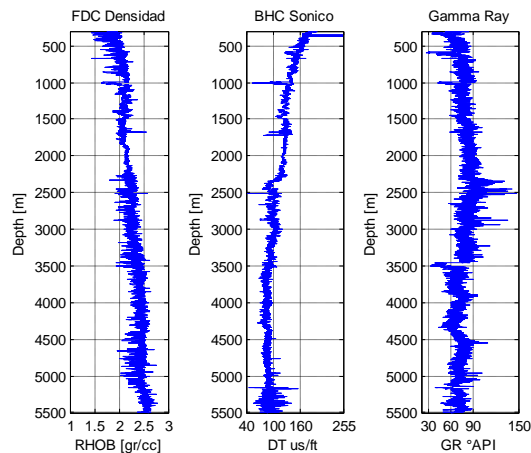
**Result and comparing with the Record Porosity curve**



Figure 4. Porosity Curve marked in Red the Trend of Minimum Values

**Well 3G – Gamma Curve and Transit Time**

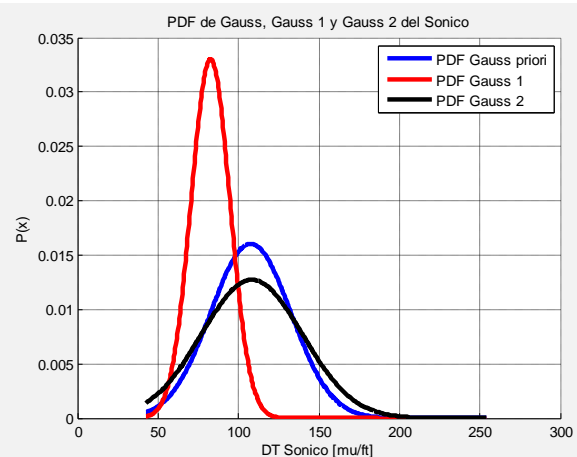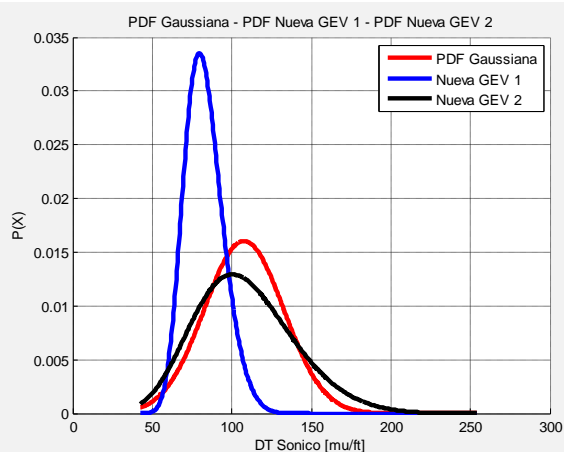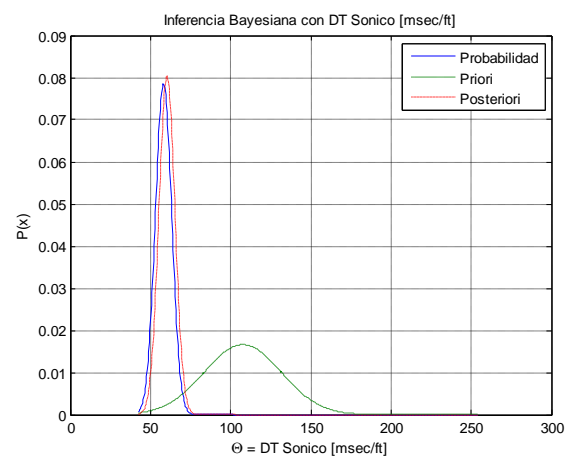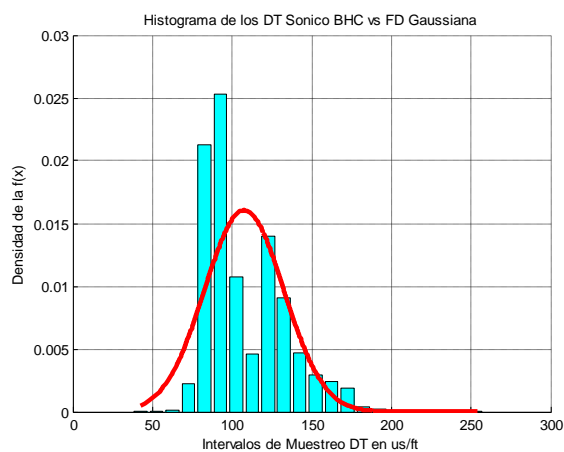Figure 5. Gamma Histogram and Gaussian Mixture marked in Red the Trend of Minimum Average Values

## With its Sonica Curve

**Results**



Figure 6. Histogram of Transit Time and Gaussian Mixture marked in Red the Trend of Average Minimum Values

**Explication**



It is observed that the PDF in blue continues the kernel of the Sand for this deposit gives greater probability of the values found in the GR and in dotted blue that of the Clay with certain trend values



It is observed that the PDF in blue continues the kernel of the Sand for this deposit gives greater probability of the values found in the Sonic and in dotted blue that of the Clay with even more increased trend values

Figure 7. Record of Transit Time and Gamma marking the Trend of Average Minimum Values found.

## Normal Transformation for Skewed Data

Sometimes one has the problem of making two samples comparable, that is, comparing the measured values of a sample with respect to its (relative) position in the distribution. (http://www.statistics4u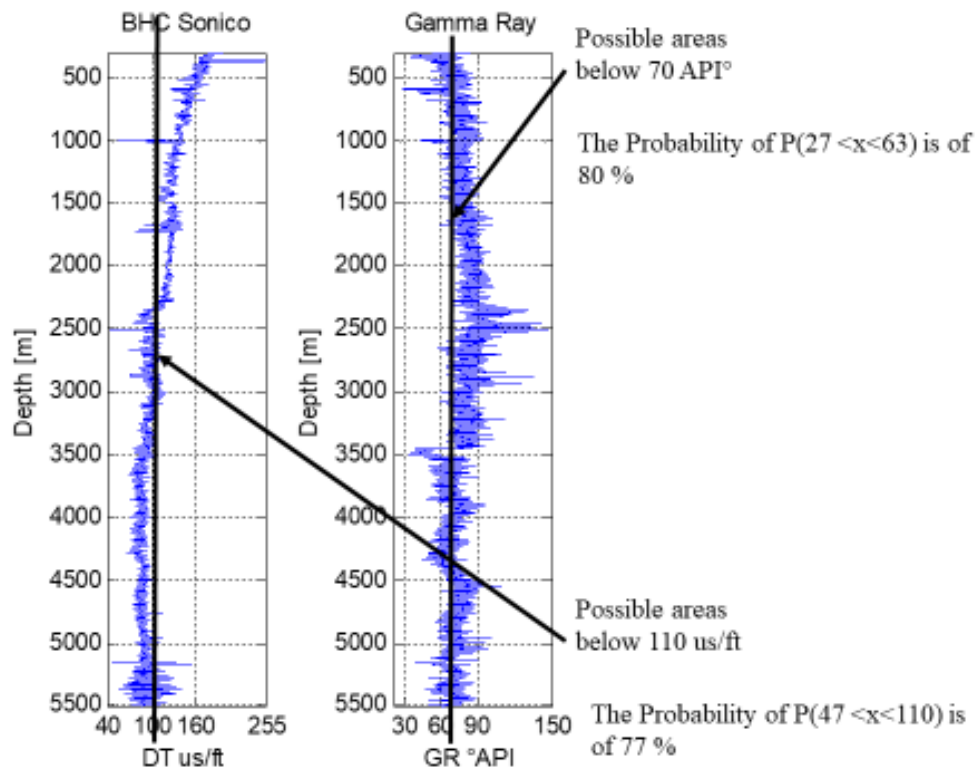.info/fundstat_eng/ee_ztransform.html) A frequently used aid is the z transform which converts the values of a sample into z scores:

With

$$Z_i = \frac{x_i - \mu}{\sigma}$$

The z transform is also called standardization or autoscaling. The z-scores are made comparable by measuring the observations in multiples of the standard deviation of that sample. The mean of a z-transformed sample is always zero. If the original distribution is normal, the z-transformed data belongs to a standard normal distribution ($\mu = 0$, s = 1).

The following example demonstrates the effect of data standardization. Suppose we have two normal distributions, one with a mean of 10.0 and a standard deviation of 30.0 (top left), the other with a mean of 200 and a standard deviation of 20.0 (top right). Standardization of both data sets results in comparable distributions since both z-transformed distributions have a mean of 0.0 and a standard deviation of 1.0

Figure 6. Histogram and Normal Transformation of Data
(http://www.statistics4u.info/fundstat_eng/ee_ztransform.html)

In some published articles, you can read that the z-scores are normally distributed. This is incorrect: the z transform does not change the shape of the distribution, it only fits the mean and standard deviation. In pictorial terms, the distribution is simply shifted along the x-axis and expanded or compressed to achieve a zero mean and a standard deviation of 1.0.



Figure 1: Procedure for transforming core porosity values, $z$, to normal score values, $y$.

Figure 7. Histogram and Normal Transformation of Data
http://www.geostatisticslessons.com/lessons/normalscore

**Table 5. Now using the case of the porosity of the Well 2 Apertura.**

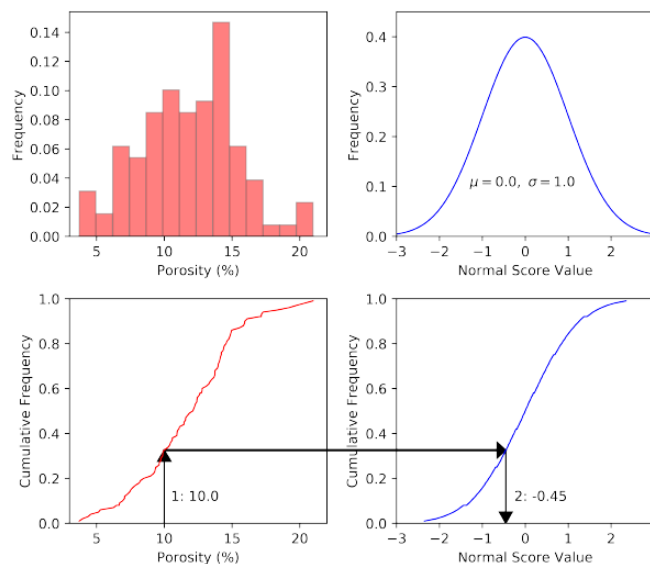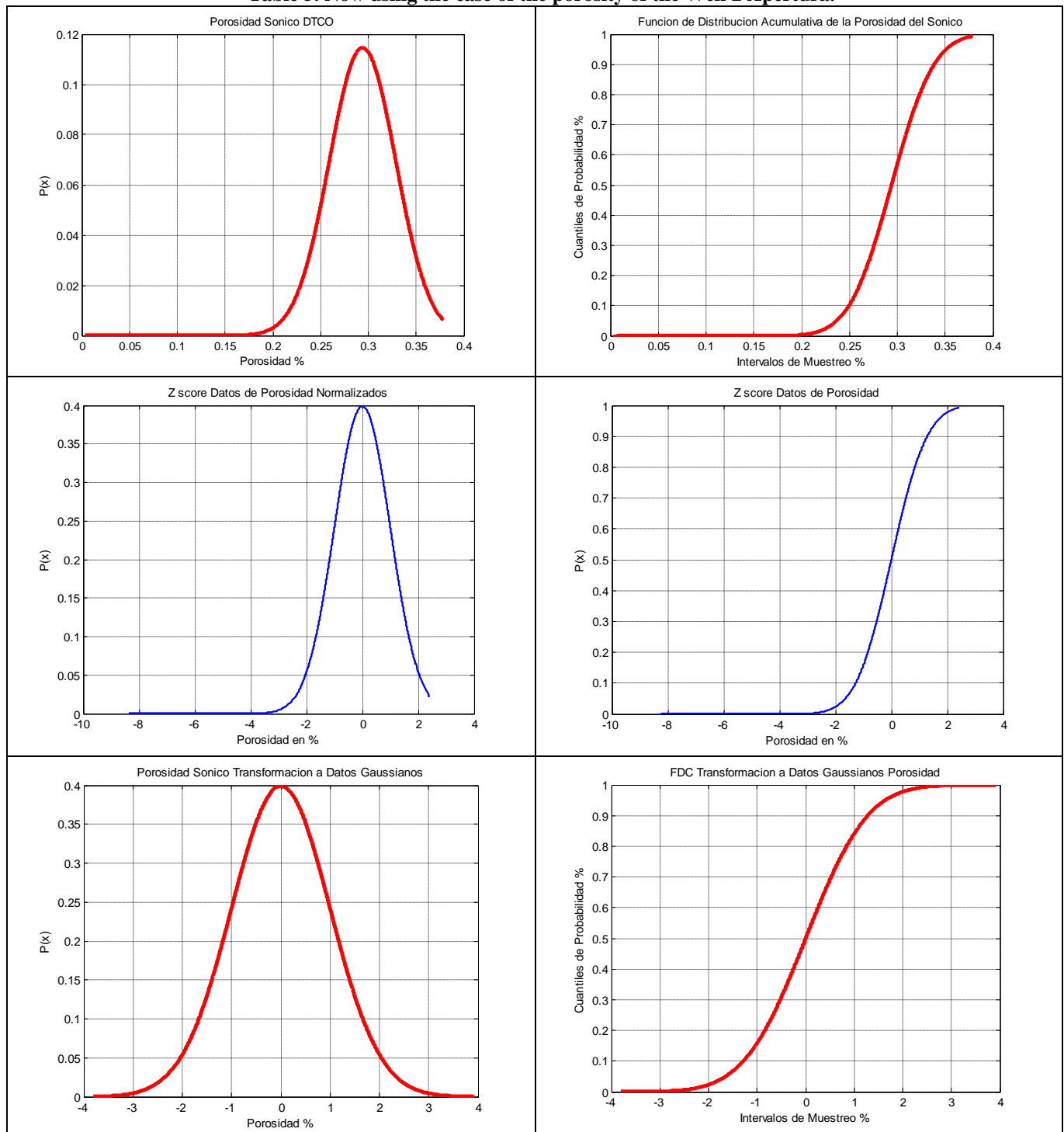It can be seen that the mean calculated from the previous adjustment is 0.30 approx. And in the transformation in both cases it is in Mean 0, the Z score of Matlab and the modified code below were used to generate the Gaussian data.

**Source code:**

```
% Data Transformation to a Gaussian
% Modified by M.C Zenteno Jimenez Jose Roberto, 02/09/2020
% Original
% Geoffrey Goldman, 6/4/2015
% https://la.mathworks.com/matlabcentral/fileexchange/50412-transform-data-
to-gaussian-distribution
clc
clear all
close all

% we generate data or a vector with the data to transform
% P = P '; Vector with data
% N = length (P); % length
N=1000;  % number of data points
data=zeros(1,N);
data=normrnd(0,1,[1,N]);% a random normal of mean is generated 0 and SD of
1

figure(1)
hist(data,50)
```

```matlab
title('Data Histogram')
xlabel('Sample Data')
ylabel('Frequencies ')
grid on

% the vector is made linear over the data
Pico=max(data);
Minimo=min(data);
lin_array=linspace(Minimo,Pico,N);
% CPDF
sdata=sort(data);  %
data_uniforme=interp1(sdata,lin_array,data,'spline');
% transformation to a uniform distribution
% with interpolation to data in splines
% Data is now equal to lin_array
% normalizing data between 0 and 1
u1_data=data_uniforme/(max(data_uniforme)-min(data_uniforme));
u1_data=u1_data -min(u1_data);

figure(2)
hist(u1_data,50,'r')
title('Histogram of the Data to a Uniform Distribution')
xlabel('Sample Data')
ylabel('Frequencies ')
grid on

% Gaussian pdf y cpdf
pro2=0.1;
lin2_array=-5:pro2:5;
%Gaussian
gauss_array=(1/(2*pi))^0.5*exp(-((lin2_array-mean(lin2_array)).^2)/2);  %
area_1=sum(gauss_array)*pro2;  %
cpdf_gauss=cumsum(gauss_array*pro2);  % cpdf
gauss_data=interp1(cpdf_gauss,lin2_array,u1_data);

figure(3)
hist(gauss_data,50,'r')
title('Histogram transformed to Gaussian')
xlabel('Sample Data')
ylabel('Frequencies ')
grid on
```

## Conclusions

According to the methodology that was now exposed the fourth case for Gaussian data and now finding the minimum values I have inferred the highest probability of those values that we are finding, we can see that it coincides with the values found as the case of Well 3 where There is a Porosity already calculated within the registration result and the comparison with the minimum Porosity values are consistent and coincide with the result shown. If you know the values that must be found within your respective reservoir, both of Porosity and Gamma values or Transit Time with this proposed method, the result can be further emphasized, with the case of Well 2 Apertura a normal case was shown in where Porosity is skewed and does not coincide from the beginning with the Gaussian Fit, therefore a normalization of the curve was made through the Z score statement and the exposed code, the original Source code and the modification are included, it was compared the two results giving a Good normalization, it must be said that if the same Bayesian Inference process is carried out with the normalization of the distribution functions obtained as the mean is zero and the Variance 1, it remains in the domain of the extreme distribution functions and their means and variances will have the effect on the skewed Gaussian and will approximate the approximate low values.

## References

[1]. A.J. Jakeman, J.A. Taylor, R.W. Simpson, Modeling distributions of air pollutant concentrations - II. Estimation of one and two parameters statistical distributions, Atmos. Environ., 20 (1986) 2435-2447.

[2]. Bayesian Online Changepoint Detection Ryan Prescott Adams, David J.C. MacKay https://arxiv.org/abs/0710.3742

[3]. Forecasting and Estimating Multiple Change-point Models with an Unknown Number of Change-points Gary Koopyand, Simon M. Potterz2006

[4]. Casella, G and Robert, C. Introducing Monte Carlo Methods with R (Use R)

[5]. Compact approximations to Bayesian predictive distributions Edward Snelson, Zoubin GhahramaniICML2005

[6]. Gumbel, E.J., 1958. Statistics of Extremes. Columbia University Press, New York, p. 164.

[7]. Hoel, P; Port, S and Stone, C. Introduction to Stochastic Processes

[8]. P.G. Georgopoulous, J.H. Seinfeld, Statistical distributions of air pollutant concentrations, Environ. Sci. Technol., 16 (1982) 401A-416A.

[9]. Roberts, E.M.,1979. Review of statistics of extreme values with applications to air quality data, part II. Applications. Journal of Air Pollution Control Association 29, 733–740.

[10]. Trabajo presentado en el Congreso de la Unión Geofísica Mexicana 2017 Pronostico de Concentraciones de Ozono por Distribuciones de Probabilidad para la CDMX https://www.raugm.org.mx/2017/pdf/constancia.php?clave=809

[11]. Prescott, P., and A. T. Walden, Maximum-likelihood estimation of the parameters of the three-parameter generalized extreme-value distribution from censored samples, J. Stat. Comput. Simul., 6, 241–250, 1983.

[12]. Otten, A., and M. A. J. Van Montfort, Maximum-likelihood estimation of the general extreme-value distribution parameters, J. Hydrol., 47, 187–192, 1980.

[13]. Zenteno Jiménez José Roberto, Prediction of Concentrations of Ozone Levels in México City using Probability Distribution Functions, International Journal of Latest Research in Engineering and Technology (IJLRET) || Volume 04 - Issue 07 || July 2018 || PP. 35-45

[14]. Zenteno Jiménez José Roberto. A Methodology for Obtaining news Probability Distributions Functions Normal and Extreme Value for Bayesian Inference and Stochastic Mixed Gaussian Case One: For Daily Concentration Data Maximum Ozone. International Journal of Latest Research in Engineering and Technology (IJLRET)|| Volume 04 - Issue 11 || November 2018 || PP. 15-35

[15]. Statistical Modeling and Computation, Dirk P. Kroese – Joshua C.C. Chan, Springer Ed. 2014, Bayesian Inference Chapter 8, 236 page.

[16]. Zenteno Jiménez José Roberto Prediction of Concentrations of Suspended Particle Levels of 2.5 micrometers (PM2.5) in Mexico City with Probability Distribution Functions and its Trend. http://www.ijlret.com/Papers/Vol-05-issue-04/1.B2019025.pdf

[17]. Heat Islands in México City: Perspective from Remote Sensing Satellite Images, Author: Fernando Mireles Arellano, Amanda Oralia Gómez González, Carlos Hernández López; International Journal of Latest Research in Engineering & Technology (IJLRET) || Volume 04 - Issue 10 || October 2018 || PP. 1-12