

## Extraction of Pectin in Seriguela (*Spondias Purpurea* L.) Using Experiments Designs: a Tutorial Using the Free Software "R"

### Extração de Pectina em Seriguela (*Spondias purpurea* L.) Usando Planejamentos Experimentais: um Tutorial Utilizando o Software Gratuito "R"

Luiz B. S. Filho,<sup>a,b</sup>  Ronaldo C. Coelho,<sup>b,\*</sup>  Tiago L. S. Coêlho,<sup>c</sup>  Edvani C. Muniz,<sup>d,e</sup>  Herbert de S. Barbosa<sup>e</sup> 

<sup>a</sup> Universidade Federal do Piauí, Departamento de Química, Grupo de Estudos em Bioanalítica – GEBIO, Zip Code 64049-550, Teresina-PI, Brazil

<sup>b</sup> Instituto Federal do Piauí, Departamento de Formação de Professores, Zip Code 64053-390 Teresina-PI, Brazil

<sup>c</sup> Instituto Federal do Amapá, Departamento de Química, Grupo de Pesquisa em Mineração, Materiais e Meio Ambiente, Zip Code 68909-398 Amapá-AP, Brazil

<sup>d</sup> Universidade Federal do Piauí, Departamento de Química, Zip Code 64049-550, Teresina-PI, Brazil

<sup>e</sup> Universidade Estadual de Maringá, Departamento de Química, Zip Code 87020-900, Maringá-PR, Brazil

E-mail: [ronald@ifpi.edu.br](mailto:ronald@ifpi.edu.br)

Submissão: 5 de Abril de 2024

Aceite: 13 de Novembro de 2024

Publicado online: 17 de Dezembro de 2024

This manuscript describes an experiment that emphasizes concepts related to the design of experiments using the computational environment R, which can be used by beginners, undergraduate students and postgraduate researchers. The pectin extraction experiment using seriguela bark (*Spondias purpurea* L.) was chosen because most of the necessary materials are available in analytical laboratories. Additionally, this tutorial provides an easy guide for installing R and other dependencies. The approach used in this tutorial introduces multivariate concepts, both in screening and in experimental optimization. First, a fractional factorial design is applied and the model is evaluated to define the experimental domain. Second, an extension of the factorial design is described, thus providing a Box-Behnken design for experimental optimization. The quadratic model is then introduced and used to construct the response surface.

**Keywords:** Free software; design of experiments; multivariate optimization; pectin.

## 1. Introduction

Chemometrics can be defined as the application of mathematical, statistical, and computational methods to investigate, interpret, classify, and predict datasets of interest.<sup>1</sup> Among the various subareas of the field, we can highlight the design of experiments (the focus of this work), a methodology that involves a multivariate approach aiming to maximize the relationship between the quality of information about a system or chemical process, that is, the main objective of designing experiments is to obtain the maximum amount of “information” by limiting the number of observations needed.<sup>2</sup>

This information is usually obtained by planning and rationally selecting experiments to obtain the best knowledge of the system.<sup>3</sup> Experiment planning can start with the study of statistically independent variables, that is, those whose values can change and be controlled independently of each other.<sup>4</sup> The selection of the type of planning to be used depends on the researcher’s objective and the project’s stage. For an initial study with the objective of screening variables and identifying those variables that have the greatest influence on the response, fractional or saturated planning should be used.<sup>5</sup> For a fine adjustment of the significant variables, a three-level factorial can be used, central composite design – CCD<sup>6</sup>, Box-Behnken.<sup>8</sup>

It should be noted that an important stage is the evaluation of the models constructed by the different types of design, that is, their adequacy for the responses obtained experimentally, will dictate their forecasting capacity.<sup>9</sup> This diagnosis can be performed in several ways, the most common being analysis of variance (ANOVA), evaluation of the residual graph (differences between the values obtained experimentally and those predicted by the model), and graphing of experimental values vs predicted values generated by the model.<sup>10</sup>

Several studies have used various designs of experiments, and an increasing number of new users from academia (undergraduate, graduate students, and researchers) and industry are using this methodology.<sup>11,12,13,14,15,16,17,18</sup> However, these new users are subject to initial limitations due to the lack of prior theoretical knowledge and ability to operate *software* and computational environments that are indispensable for their application. In this context, it is important to understand the calculations performed by the *software*, which is fundamental for evaluating the results obtained, as well as for questioning how such *software* performs them.

Many *software* packages that can be used in routine handling experiments are available on

the market, including MATLAB, Design Expert, Statistica, Minitab, Pirouette, and Unscrambler, as well as free *software* such as Octave and R.<sup>12</sup> Although the creation of tutorials involving the design of experiments is not something new, some can be found in the scientific literature involving *software* such as MATLAB and Octave<sup>19–21</sup>, however, up-to-date and comprehensive tutorials involving designing experiments in open-source R *software* were not found in our searches.

In this way, the present work aims to provide a tutorial for the design of experiments using the R environment, in addition to showing the basic commands step by step in a direct and practical way to perform all the calculations involved.

## 2. Experimental

### 2.1. Sample preparation

Fresh seriguela (*Spondias purpurea* L.) were purchased at the fruit market in the city of Teresina-PI, Brazil. The samples were peeled, and the obtained peels were subsequently dried in an oven at 45°C for 96 h.<sup>14</sup> After drying, the peels were ground using a blender and stored in the dark and in a dry place.

### 2.2. Ultrasound assisted extraction

Approximately 0.5 g of sample was mixed with a citric acid solution (Merck Chemical Co., Darmstadt, Germany) in a 50 mL polypropylene tube and then placed in an ultrasonic bath (Ultrasonics Cleaner – Soni-Tech®) using the experimental design conditions (Table 1S).<sup>14</sup>

The acid extract was cooled to 4°C for approximately 2 h. Then, the mixture was centrifuged at 3500 rpm (1090 xg) for 10 minutes. To the pectin-containing supernatant, ethyl alcohol (95%) was added at a 1:3 (v/v) ratio (one part of the pectin containing solution and three parts alcohol). After 2 hours of rest, a pectin precipitate was obtained, which was centrifuged at 3500 rpm (1090 xg) for 10 minutes. Finally, the supernatant was discarded, and the resulting material was dried in an oven at 50°C until constant weight.

The pectin extraction efficiency was calculated using Eq. 1, where  $Y$  is the extracted pectin yield as a percentage (%),  $mf$  is the amount of extracted pectin in g and  $mi$  is the initial amount of ground seriguela peel.

$$Y(\%) = \frac{mf}{mi} \times 100 \quad (1)$$

### 2.3. Design of experiments

A 2<sup>5-1</sup> fractional factorial design (Table 1S) was applied to screen the variables independent (hydrogenation potential

(pH), extraction temperature (T), extraction time (t), ratio between the mass of ground seriguela shell and water (s/l), and potency of the ultrasonic bath (Pw)) chosen in view of studies reported in the literature.<sup>13,14</sup> In the fractional design, the extraction yield (y(%)) was used as a response (dependent variable), and the generatrix used was 12345, where the fifth variable was obtained by multiplying variables 1, 2, 3 and 4 (2<sup>5-1</sup>). All the experiments were performed randomly in order to minimize the effect of unexplained variability in the observed responses due to systematic errors.<sup>22</sup>

The variables were coded according to Eq. 2, where  $x$  is the coded value,  $X_i$  is the corresponding real value,  $X_0$  is the real value at the central point and  $\Delta X$  is the increment of  $X_i$  corresponding to a variation of 1 unit of  $x$ .

$$x = \frac{X_i - X_0}{\Delta X} \quad (2)$$

After the screening step, only the significant variables were selected, and new levels were defined for the refinement of the extraction method using a Box-Behnken design (BBD). The statistical model obtained was validated by analysis of variance (ANOVA), analysis of residual graphs, and experimental vs predicted values. The optimal conditions for each variable were evaluated using the response surface methodology.

## 3. Results and Discussion

To process the results obtained, it is necessary to follow the installation guide for the R, RStudio, and Rtools extension programs indicated in the supplementary material. It is worth mentioning that it is necessary to install and load all the packages that were used to execute all the functions and commands in this tutorial. For this purpose, after starting RStudio, clicking on “Ctrl+Shift+N” (shortcut informed in the item – RStudio Interface of the Installation Guide), the Editor window will appear.

In the editor window, according to Table 1, type the argument to install the “FrF2” package, is typed, the cursor is positioned inside the argument, and the “Ctrl+enter” is pressed. For the package to perform its functions, it is necessary to load it. These same steps were followed for the other packages shown in Table 1.

**Table 1.** Codes for installing and loading packages.

	Code
<b>Install FrF2</b>	install.packages(“FrF2”, dependencies = TRUE)
<b>Load FrF2</b>	library(FrF2)
<b>Install rsm</b>	install.packages(“rsm”, dependencies = TRUE)
<b>Load rsm</b>	library(rsm)
<b>Install ggpubr</b>	install.packages(“ggpubr”, dependencies = TRUE)
<b>Load ggpubr</b>	library(ggpubr)

More details on the composition of all arguments in the codes for this tutorial can be found in the supplementary material. It is recommended that the organization of the experimental data obtained be carried out in a spreadsheet in .csv or .xlsx formats, in addition to files in .txt format, which are saved in a folder on your computer to facilitate manipulation. Alternatively, data organization can be performed in RStudio itself.

Now, the command with the “Ctrl+Shift+H” keys is used as a shortcut, and the working directory is changed to the same folder, in which the files are saved in .csv, .xlsx, or .txt format. This makes it easy to import these files for manipulation in RStudio.

### 3.1. Fractional factorial design $2^{5-1}$

To start processing the data in the *software*, it is necessary to import the planning matrix similar to Table 1S (Supplementary Material). In the editor window in RStudio, always type what is highlighted in blue and follow the same information already mentioned for the arguments in Table 1.  
# Import:

```
M1 <- read.csv("triagPec.csv",header = F,dec = ".",sep = ";")
```

The “Ctrl+enter” keys were pressed, and the data were imported into the .csv format. M1 was the name chosen for the object created from the source file of the matrix, “triagPec.csv”. Note that this created object will appear in the global environment window. Then, the user clicks on the object created to view it directly in the editor window.

As shown in Table 1S, the influence of 5 factors (pH, temperature, time, potency, and ratio) on the extraction of pectin from the bark of the seriguelas was evaluated. In this section, a  $2^{5-1}$  design was used to study the influence of the 5 variables on the pectin extraction yield (Table 1S). In all, 19 experiments were performed with three replicates at the central point to determine the experimental error.

Note in Table 1S that decimal numbers in R must be separated by a period (.) and not a comma (,).

#### 3.1.1. Creating a matrix from the “FrF2” package

The FrF2 package can be used to analyze data from a two-level fractional factorial design. In addition to evaluating the effects of the experiment, interaction graphs and main effects are constructed.

#### 3.1.2. Planning without center points (PwCP)

First, planning without central points will be carried out. In this case, it will help later on to plot the normal probability graph (Figure 4), since the function to be executed does not allow planning with central points to be used. Type and execute the command in blue and then press “Ctrl+enter”:  
# PwCP:

```
frac1 <- FrF2(nruns = 16, nfactores = 5, factor.names = c("x1", "x2", "x3", "x4", "x5"), randomize = F, alias.info = 3)
```

Note that the “frac1” object created will appear in the global environment window. Then, the user clicks on the object created to view it from a tab in the editor window. The generic function of the editor is to summarize the planning structure without central points.

```
# Planning Summary:  
summary(frac1)
```

Click Ctrl + enter and note that the results will be shown directly in the console and will not be saved in the global environment, as shown in Figure 1.

#### 3.1.3. Creating the response vector

This vector must contain the responses in the same order as the experiments performed or contained in the “M1” matrix. There are two ways to create the response vector, the first is by writing the arguments with the experimental responses, and the second is by indexing these responses from the M1 matrix. Here, in this work, the second option was used. To do this, enter the arguments:

```
# Response Vector:  
y1 <- M1[1:16,6]
```

Then, press “Ctrl+enter”. Note that the “y1” object created will appear in the global environment window.

#### 3.1.4. Regression model

To create the planning regression model without central points, the script was written in blue, and the “Ctrl+enter” command was subsequently used. To view the regression summary (Figure 2), the same procedure was repeated for each PwCP.

```
# Regression model:  
lm1 <- lm(y1 ~ .^2, data = frac1)  
#Summary of Regression data:  
summary(lm1)
```

As shown in Figure 2, the experimental error and significance cannot be calculated, because genuine replicates and central points were not used.

#### 3.1.5. Analysis of interactions

To view the graphs of the main effects and interactions, type and then press “Ctrl+enter” for each argument below:  
# Main effects:

```
MEPlot(lm1, main = “Gráfico dos Efeitos Principais”,  
pch = 16, cex.main = 1.3, lwd = 2.5, cex.xax = 1.7, mgp.  
ylab = 4, cex.yax = 1.7)
```

```
# Effect of interactions
```

```
IAPlot(lm1, main = “Gráfico das Interações”, pch = c(25,16),  
cex.main = 1.8, lwd = 2, cex.xax = 1.7, cex.lab = 1.7,  
cex = 1.5)
```

It is noteworthy that the “MEPlot” and “IAPlot” functions of the FrF2 package accept metrics only from factorial planning without central points. The figures were

```

Call:
FrF2(nruns = 16, nfactores = 5, factor.names = c("x1", "x2",
"x3", "x4", "x5"), randomize = F, alias.info = 3)

Experimental design of type FrF2
16 runs

Factor settings (scale ends):
  x1 x2 x3 x4 x5
1 -1 -1 -1 -1 -1
2  1  1  1  1  1

Design generating information:
$legend
[1] A=x1 B=x2 C=x3 D=x4 E=x5

$generators
[1] E=ABCD

Alias structure:
$F12
[1] AB=CDE AC=BDE AD=BCE AE=BCD BC=ADE BD=ACE BE=ACD CD=ABE CE=ABD DE=ABC

The design itself:
  x1 x2 x3 x4 x5
1 -1 -1 -1 -1  1
2  1 -1 -1 -1 -1
3 -1  1 -1 -1 -1
4  1  1 -1 -1  1
5 -1 -1  1 -1 -1
6  1 -1  1 -1  1
7 -1  1  1 -1  1
8  1  1  1 -1 -1
9 -1 -1 -1  1 -1
10  1 -1 -1  1  1
11 -1  1 -1  1  1
12  1  1 -1  1 -1
13 -1 -1  1  1  1
14  1 -1  1  1 -1
15 -1  1  1  1 -1
16  1  1  1  1  1
class=design, type= FrF2

```

Figure 1. Structure Summary of the PwCP

```

Call:
lm.default(formula = y1 ~ x1 + x2 + x3 + x4 + x5 + x1 * x2 +
  x1 * x3 + x1 * x4 + x1 * x5 + x2 * x3 + x2 * x4 + x2 * x5 +
  x3 * x4 + x3 * x5 + x4 * x5, data = frac1)

Residuals:
ALL 16 residuals are 0: no residual degrees of freedom!

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  25.8062          NA      NA      NA
x11           3.6163          NA      NA      NA
x21           0.8313          NA      NA      NA
x31           2.2838          NA      NA      NA
x41           2.6575          NA      NA      NA
x51           0.3625          NA      NA      NA
x11:x21       0.2687          NA      NA      NA
x11:x31      -2.9613          NA      NA      NA
x11:x41       3.0325          NA      NA      NA
x11:x51      -0.3850          NA      NA      NA
x21:x31       0.1987          NA      NA      NA
x21:x41      -0.1750          NA      NA      NA
x21:x51      -0.2225          NA      NA      NA
x31:x41      -0.1275          NA      NA      NA
x31:x51       0.6675          NA      NA      NA
x41:x51       0.1063          NA      NA      NA

Residual standard error: NaN on 0 degrees of freedom
Multiple R-squared:  1,    Adjusted R-squared:  NaN
F-statistic:  NaN on 15 and 0 DF,  p-value: NA

```

Figure 2. Summary of the significance of the estimated coefficients

automatically plotted in the output window (output). The main effects demonstrate the individual impact of each factor on pectin yield (see Figure 34S(a)).

When analyzing each variable, it was observed that there was a positive impact when the pH, temperature,

time and ultrasound power moved to the upper level of the factor, remaining practically unchanged with the variation in the ratio (variable x5). The interaction matrix given in Figure 34S(b) indicates that the interactions pH: time and pH: power present a greater relationship by crossing the

curves under higher and lower levels of factors.

It is noteworthy that the results presented thus far clearly demonstrate that the order of importance of the factors are pH ( $x_1$ ) > ultrasound power ( $x_4$ ) > time ( $x_3$ ) > temperature ( $x_2$ ) for the planning carried out, since the factor ( $s/l$ ) was not significant at the 95% confidence level.

### 3.1.6. Planning with center points (PWCP)

The application of planning with central points translates into an advantage in being able to calculate the experimental error and the significance of the parameters. To view it, type the following argument and then press “Ctrl+enter”.

# PWCP:

```
Frac2 <- FrF2(nruns = 16, ncenter = 3, nfactors = 5, XXXator.
names = c("x1", "x2", "x3", "x4", "x5"), randomize = F,
alias.info = 3)
```

Note that the difference between the script for planning without a center point was the inclusion of information indicating the number of center points now used ( $n_{center} = 3$ ). Additionally, the object “frac2” created will appear in the window of the global environment. Then, the user clicks on the object created to view it from a tab in the editor window. To view the results summary, type the following function in the editor and press Ctrl + enter.

# Planning Summary:

```
summary(frac2)
```

Note that the planning structure results will be shown directly in the console and will not be saved in the global environment. Figure 35S is a summary of the levels, and associated factors and presents the structure of the confounding patterns (contrasts) along with the coded planning matrix.

Now, to insert the answers, as was done in planning without a central point, type the following script in the function editor and then press “Ctrl+enter”.

# Response Vector:

```
y2 <- M1[,6]
```

The “y2” object created will appear in the global environment window.

### 3.1.7. Calculation of the main and interaction coefficients

With the experiments described in Table 1S, it is possible to calculate 15 effects - 5 main effects – one for each variable individually ( $x_1, x_2, x_3, x_4$  and  $x_5$ ) and 10 secondary interactions ( $x_1x_2, x_1x_3, x_1x_4, x_1x_5, x_2x_3, x_2x_4, x_2x_5, x_3x_4, x_3x_5$  and  $x_4x_5$ ). Since the values of the effects of tertiary (123, 124, 125, 134, 135, 145, 234, 235, 245, 345) and quaternary (1234, 1235, 1245, 1345, 2345) interactions are not defined, as they are confused with secondary interactions and main effects, respectively (see Figure 35S), as it is a  $2^{5-1}$  fractional factorial design.

In the editor window, type the following command is used to determine the linear regression model and thereby determine the coefficients and effects of the design with center points. Always remember to press “Ctrl+enter” after writing the argument.

# Regression model:

```
lm2 <- lm(y2 ~ .^2, data = frac2)
```

Since third- and fourth-order interactions are confounded with second-order and main effects, respectively, there is no need to obtain the complete model. The “lm2” object created is a list that will appear in the global environment window. Then, to view it directly in the editor window, the experimenter clicks on the created object. A summary of the results can be viewed by typing the script and clicking ctrl+enter.

# Summary of Regression data:

```
summary(lm2)
```

The results summary should appear directly on the console similar to that shown in Figure 3.

```
Call:
lm.default(formula = y2 ~ x1 + x2 + x3 + x4 + x5 + x1 * x2 +
  x1 * x3 + x1 * x4 + x1 * x5 + x2 * x3 + x2 * x4 + x2 * x5 +
  x3 * x4 + x3 * x5 + x4 * x5, data = frac2)

Residuals:
    1     2     3     4     5     6     7     8     9    10    11    12    13    14    15
0.0952 0.0952 0.0952 0.0952 0.0952 0.0952 0.0952 0.0952 0.0952 0.0952 0.0952 0.0952 0.0952 0.0952 0.0952
    16    17    18    19
0.0952 -0.2811 -0.4711 -0.7711

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  25.7111      0.1351  190.308 3.2e-07 ***
x1           3.6163      0.1472   24.563 0.000148 ***
x2           0.8312      0.1472    5.646 0.010996 *
x3           2.2837      0.1472   15.512 0.000582 ***
x4           2.6575      0.1472   18.051 0.000371 ***
x5           0.3625      0.1472    2.462 0.090688 .
x1:x2        0.2688      0.1472    1.825 0.165420
x1:x3       -2.9612      0.1472  -20.114 0.000269 ***
x1:x4        3.0325      0.1472   20.598 0.000250 ***
x1:x5       -0.3850      0.1472   -2.615 0.079339 .
x2:x3        0.1987      0.1472    1.350 0.269845
x2:x4       -0.1750      0.1472   -1.189 0.320091
x2:x5       -0.2225      0.1472   -1.511 0.227890
x3:x4       -0.1275      0.1472   -0.866 0.450186
x3:x5        0.6675      0.1472    4.534 0.020081 *
x4:x5        0.1062      0.1472    0.722 0.522650
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.5889 on 3 degrees of freedom
Multiple R-squared:  0.9986,    Adjusted R-squared:  0.9913
F-statistic: 138.3 on 15 and 3 DF,  p-value: 0.0008847
```

Figure 3. Summary of the significance of the estimated coefficients

According to the summary table of the values of the coefficients and their significance, (see Figure 3), the factor “ $x_5$ ” (ratio) was not significant at the 95% confidence level ( $p < 0.05$ ).

To obtain the effects significance graph, it is necessary to create a new plan from the object “frac2” and the answer “y2”. Type in the editor and then press “Ctrl+enter”:

```
# New planning:
plan1 <- add.response(design = frac2, response = y2)
```

With the new plan created, to plot the effects significance graph from the “plan1” object, type in the editor and press “Ctrl+enter”:

```
#Plot the graph:
DanielPlot(plan1, code = FALSE, half = TRUE, alpha = 0.05, pch = 18, main = “Gráfico normal de probabilidade”, cex.main = 1.3, cex.fac = 1.2, font.axis = 2, cex.lab = 1.25, cex.pch = 1.2, font.lab = 4)
```

The figure showing the significance of the effects will appear in the “plots” tab in the output window (output). The generated graph will be similar to the one shown in Figure 4.

With the normal graph of the probability of significance of the effects (Figure 4), it is possible to notice that the most important contrasts are the main effects  $x_1$ ,  $x_2$ ,  $x_3$  and  $x_4$ , as well as the most meaningful interactions  $x_1x_3$  and  $x_1x_4$ . Nevertheless, to determine the significance of the effects, a Pareto chart can be constructed for the standardization of the effects. The following arguments are used to construct the Pareto chart or diagram and other important parameters.

```
# t critical:
t_critical <- qt(0.025, df.residual(lm2), lower.tail = F)
# Experimental variance:
MSE <- deviance(lm2)/df.residual(lm2)
# Standard error of coefficients:
SE_coef <- sqrt(MSE/16)
# t calculated:
t0 <- lm2$coefficients/SE_coef
# data frame:
t_0 <- data.frame(names(coef(lm2)),abs(t0))
```

```
# Rename table columns:
colnames(t_0) <- c(“term”, “t0”)
```

To construct the Pareto chart, type the following arguments were used to press “Ctrl+enter”:

```
# Chart creation:
pPar <- ggbarplot(data = t_0[-1,], x = “term”, y = “t0”, col = “darkblue”, fill = “lightgreen”, rotate = T, sort.val = “asc”) + theme_bw() + geom_hline(yintercept = t_critical, col = “red”)
# Visualization:
ggpar(pPar, main = “Pareto chart”, font.main = c(14, “bold”, “Black”), font.x = c(14, “bold”, “Black”), font.y = c(14, “bold”, “Black”), font.tickslab = c(14, “bold”, “Black”), xtickslab.rt = 0, ytickslab.rt = 0)
```

Figure 36S shows the Pareto chart of standardized effects corroborating the results presented in Figure 3 and Figure 4. The red line represents the calculated  $t$ .

### 3.2. Optimizing the extraction process

After sorting the variables, the system was refined by analyzing only the significant variables. For this purpose, a Box–Behnken design was applied using previous information acquired in the previous planning that made it possible to define a new experimental domain, with levels and an appropriate provisional model.

The Box-Behnken model contains  $N = (2n(n - 1)) + P_c$  experiments, where  $N$  is the number of trials to be carried out,  $n$  is the number of factors studied, and  $P_c$  is the number of central points.<sup>23</sup> For the example, there are 4 factors and 5 central points, for a total of 29 experiments.

#### 3.2.1. Importing from directory

However, the optimization will be performed at another time. When opening the RStudio program, the necessary packages must be loaded. The packages to be loaded are the same as those installed in Table 1, now, only the loads will be applied, and no new installation is needed. To change the

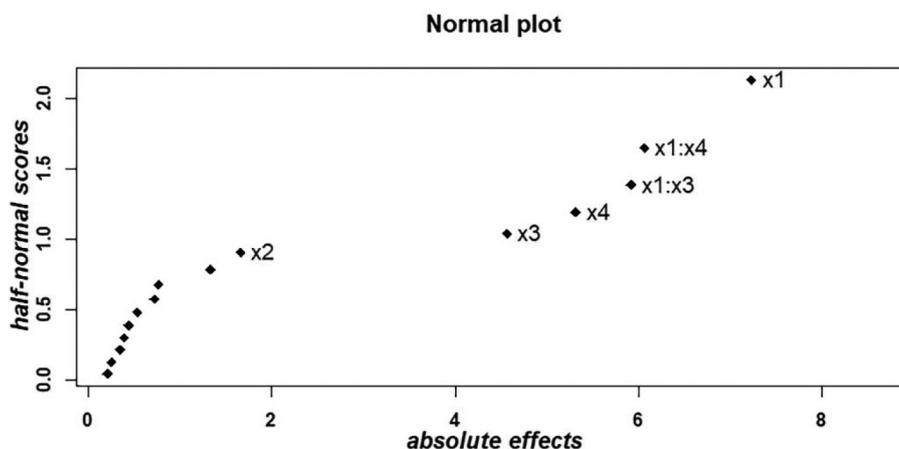


Figure 4. Normal chart of the probability of significance of effects

directory, as previously mentioned, the shortcut was typed with the “Ctrl+Shift+H” keys in the editor, and the folder that contains the files was chosen in .csv, .xlsx, or .txt format. After these steps, the planning matrix is imported; for this, the following argument is typed into the editor:

```
# Import:
M2 <- read.table("bbdPect.txt", header = F, dec = ".")
```

Then, the “Ctrl+enter” keys were pressed to finish the import. M2 is the name given to the object created from the source file of the matrix, that is, “bbdPec.txt”. Note that this created object will appear in the global environment window. Then, the user clicks on the object created to view it directly in the editor window. A table similar to Table 2S should appear.

Now, it is necessary to create the planning matrix from the “rsm” package. To do this, in the editor window, type the following argument and press “Ctrl+enter”:

```
# Creating matrix:
bbd1 <- bbd(4,5,randomize = FALSE, block = F, coding =
list(x1 ~ (pH - 2)/0.35, x2 ~ (te - 70)/5, x3 ~ (tp - 25)/10,
x4 ~ (pW - 2)/1.5))
```

The object “bbd1” created will appear in the window of the global environment. Then, the user clicks on the object created to view it directly in the editor window. Note that all factors are coded. If you want to view the created schedule directly on the console, type, and press “Ctrl+enter”:

```
# View in console:
bbd1
```

Note that the factors are now in their actual units.

### 3.2.2. Creating the Response Vector

After creating the matrix, the next step is to create the response vector. This vector must contain the responses in the same order as the experiments performed or contained in the “M2” matrix. In the function editor, type the following script and then press “Ctrl+enter”:

```
# Response Vector:
rend <- M2[,5]
```

Note that the “rend” object created will appear in the global environment window. With the created response vector, the “rend” response to the planning “bbd1” is added. This is done by typing the script below in the function editor and clicking “ctrl+enter” right after.

```
# Adding responses to bbd1:
bbd1$rend <- rend
# View in console:
bbd1
```

With the information entered, the model can be evaluated quickly. To determine the ANOVA and the coefficients of the complete model, type in the function editor below and then press “Ctrl+enter”:

```
# ANOVA and coefficients:
```

```
rsm1 <- rsm(rend ~ FO(x1, x2, x3, x4) + PQ(x1, x2, x3, x4)
+ TWI(x1, x2, x3, x4), data = bbd1)
```

The created object “rsm1” is a list that appears in the window of the global environment. Then, the created object was clicked on to visualize its constitution in the editor window. However, for the editor type, the generic function of summarizing the results is as follows:

```
#Summary of results:
summary(rsm1)
```

Note that the results will be shown directly in the console. Figure 5 presents the ANOVA and the estimated coefficients of the complete model that must be presented after writing the argument of the summary of the results.

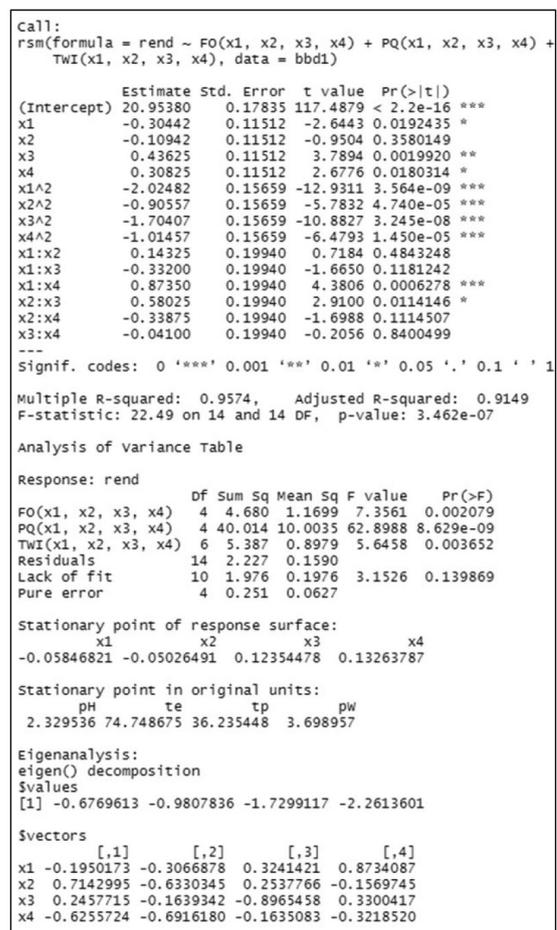


Figure 5. ANOVA summary table and regression coefficients

According to the results of the analysis of variance (ANOVA) shown in Figure 5, the high F value obtained (79.9007) and the low p-value (< 0.00056) justify that the proposed model is significant.<sup>24</sup> Furthermore, all the extraction factors significantly affected the pectin extraction yield. The coefficient of determination ( $R^2 = 0.9574$ ) implies that 95.74% of all the variations can be explained by the proposed model.<sup>25</sup> The p-value of the lack of fit test was insignificant (0.1399), which indicates that the model was

well-adjusted for the relationship between variables and the response.<sup>25</sup>

For a well-fitted model, the predicted values are expected to linearly agree with the experimental values.<sup>25,26</sup> With the coefficients presented in Figure 5, it is possible to assemble an adjusted second-order polynomial model for optimizing pectin extraction. The equation of the model becomes:

$$\text{Rend} = +20.953 - 0.304x_1 - 0.109x_2 + 0.436x_3 + 0.308x_4 + 0.143x_1x_2 - 0.332x_1x_3 + 0.873x_1x_4 + 0.580x_2x_3 - 0.338x_2x_4 - 0.041x_3x_4 - 2.024x_1^2 - 0.905x_2^2 - 1.704x_3^2 - 1.014x_4^2 \quad (3)$$

The mathematical model, Eq. 3, can be used to calculate the expected response for each of the test conditions. If the model used contains all the terms needed to adequately predict the response, the difference between the predicted and experimental values (error) should show the following behavior:

- i) must tend to a normal distribution;
- ii) should not vary depending on the expected response;
- iii) There should be no correlation with the independent variables or temporal sequence of the tests. That is, the model residuals must have the same properties as the experimental error.

To validate these assumptions, several graphs are created (residual normal probability graph, residual graph versus testing sequence, residual graph versus predicted response, and residual versus independent variables graph) that help in the validation analysis. To create the residual graph, type the following argument and then press “Ctrl+enter”.

```
# Residual chart:
```

```
plot(fitted(rsm1), resid(rsm1), abline(h = 0, lty = 2, col = "blue", lwd = 2), font.axis = 2, xlab = "Valores Previstos", cex.main = 1.3, font.lab = 4, ylab = "Resíduos", main = "Resíduos x Valores Previstos")
```

```
# Red color:
```

```
points(fitted(rsm1), resid(rsm1), pch = 21, bg = "red")
```

The plot of the residual values versus the experimentally predicted values of the generated pectin yield is similar

to that shown in Figure 37S and reveals that the errors are random, again indicating that the built model exhibits adequate linear behavior.

Now, in the graph of predicted values vs. experimental values shown in Figure 6, a linear relationship is observed, indicating that the observed values are close to the values predicted by the model. To visualize and apply the necessary arguments for its construction.

```
#Predicted x Experimental Graph:
```

```
plot(rend, rsm1$fitted.values, font.lab = 4, abline(0, 1, col = "blue", lty = 2, lwd = 2), xlab = "Experimental", ylab = "Previstos", cex.main = 1.3, main = "Experimental x Previstos", font.axis = 2)
```

```
# Red color:
```

```
points(rend, rsm1$fitted.values, pch = 21, bg = "red")
```

Now, the contour plot will be created. The arguments used in the editor are presented in Table 2, and for each written script, the press “Ctrl+enter”.

The contour graphs generated are shown in Figure 7 and Figure 38S, which show the response lines as a function of the levels of two factors.<sup>27</sup> By analyzing the y- and x-axes, we can draw straight lines and reach the regions with the highest pectin extraction efficiencies, and these regions reach the respective levels of each variable studied.

For the construction and formatting of surface graphics, the same procedure was used for contour graphics, except for the arguments in Table 3.

The surface graphics were similar to the images in Figure 8 and Figure 39S. These 3D graphs represent the response in a third dimension.<sup>27</sup> The best or optimal conditions are often derived from such graphs. However, one should be aware that if three or more factors are considered, the graphs in Figure 8 and Figure 39S represent only an (occasionally very small) part of the entire response surface in the examined domain.

The response surfaces obtained from the variables shown in Figure 8 (A and B) reveal that the increase in extraction efficiency tends toward the central point of all the variables; that is, the maximum of the curve is located at X, Y, Z, and W for the temperature, pH, time and power variables, respectively.

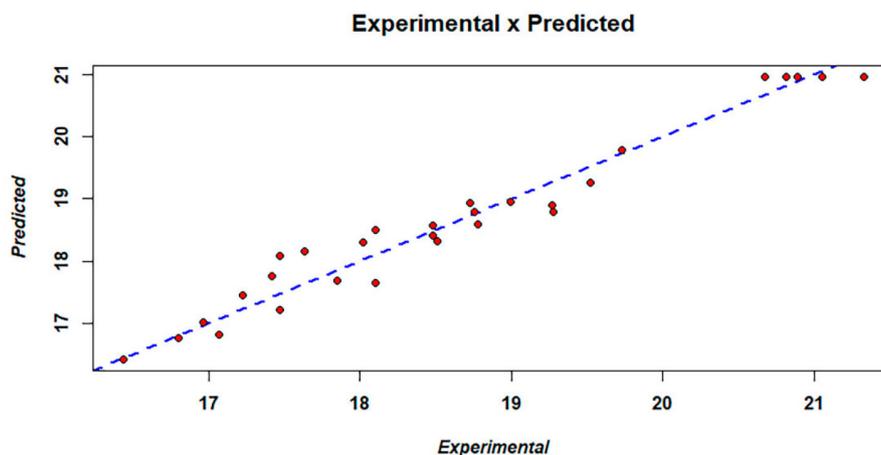


Figure 6. Plot of the predicted values vs. experimental values for optimizing pectin extraction

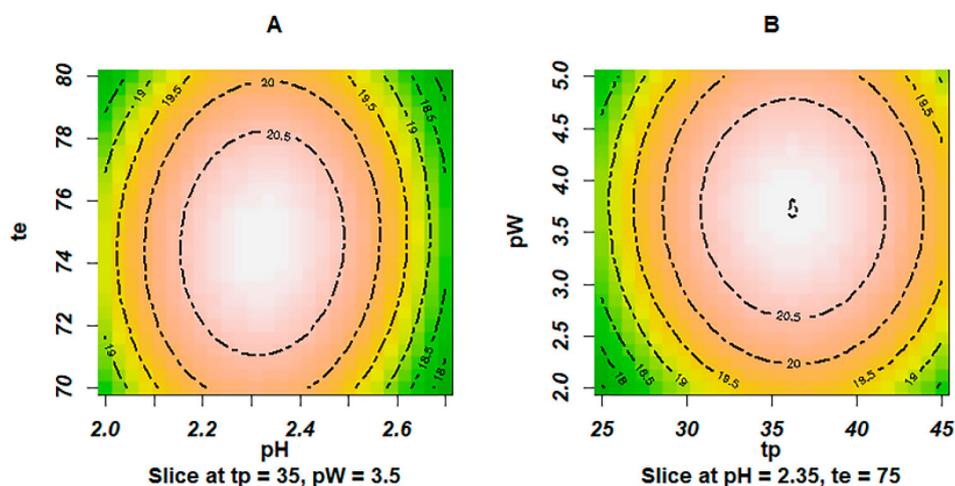


Figure 7. Contour Charts A and B

Table 2. Arguments for the construction of contour plots.

Graphical	Commands	Scripts
-	#two columns of charts	par1 <- par(mfrow = c(1,2))
A	#Command	contour(rsm1, ~x1 + x2, image = TRUE, lty = 6, cex.lab = 1.1, cex.axis = 1.1, lwd = 2, font.axis = 4, font.lab = 2)
	#Name	title(main = "A")
B	#Command	contour(rsm1, ~x1 + x3, image = TRUE, lty = 6, cex.lab = 1.1, cex.axis = 1.1, lwd = 2, font.axis = 4, font.lab = 2)
	#Name	title(main = "B")
C	#Command	contour(rsm1, ~x1 + x4, image = TRUE, lty = 6, cex.lab = 1.1, cex.axis = 1.1, lwd = 2, font.axis = 4, font.lab = 2)
	#Name	title(main = "C")
D	#Command	contour(rsm1, ~x2 + x3, image = TRUE, lty = 6, cex.lab = 1.1, cex.axis = 1.1, lwd = 2, font.axis = 4, font.lab = 2)
	#Name	title(main = "D")
E	#Command	contour(rsm1, ~x2 + x4, image = TRUE, lty = 6, cex.lab = 1.1, cex.axis = 1.1, lwd = 2, font.axis = 4, font.lab = 2)
	#Name	title(main = "E")
F	#Command	contour(rsm1, ~x3 + x4, image = TRUE, lty = 6, cex.lab = 1.1, cex.axis = 1.1, lwd = 2, font.axis = 4, font.lab = 2)
	#Name	title(main = "F")

Pectin was extracted from grape pomace using citric acid as the extracting agent, and a Box–Behnken design was used to optimize the extraction parameters to obtain a high yield of pectins with a high mean molecular weight (MW) and a high degree of esterification (DE).<sup>14</sup> The pectin and phenolic compounds were extracted from mango peels by ultrasonication. The extraction aided by the ultrasound method helped to increase the pectin yield by 50% without affecting its quality.<sup>28</sup> Pectin was extracted from tomato processing residues using ultrasound-assisted extraction and other techniques. The pectin yields were 9.30% and 25.42% for microwave-assisted and ultrasound-assisted extraction, respectively.<sup>29</sup> Notably, until now, no studies have been found in the literature on the yield of pectin extracted from seriguela (*Spondias purpurea* L.), but the yields were comparable to those previously reported for other extraction sources.

All the results, such as the tables and graphs shown in R, were compared with the figures and metrics presented in the *software* Minitab® version 17.3.1 and Design Expert® version 7.0.0 according to the supplementary data.

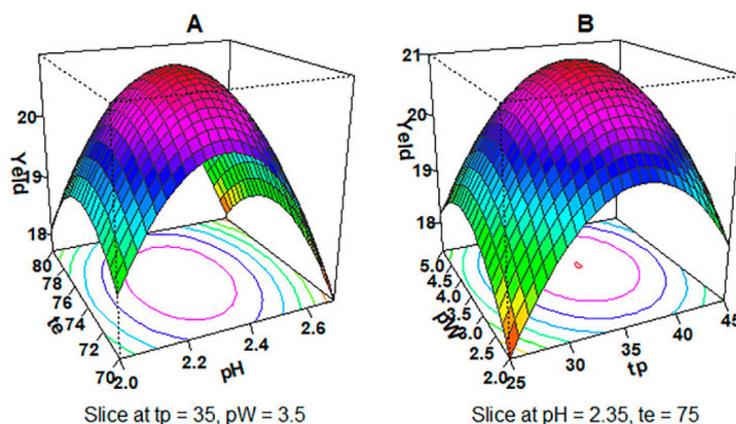
### 3.3. Stationary point and optimization via steepest ascent (maximum slope path)

To confirm the information displayed on the response surfaces and contour plot, we can use optimization by steepest ascent.

A summary (Figure 5) of the second-order model provided the results of a canonical analysis of the surface.<sup>30</sup> The analysis indicated the stationary point of the fitted surface<sup>25</sup> (pH = -0.058; temperature = -0.050; time = 0.123; and ultrasound power = 0.132) in coded units within the

**Table 3.** Arguments for the construction of surface graphics

Graphical	Commands	Scripts
-	<b>#two columns of charts</b>	<code>par2 &lt;- par(mfrow = c(1,2))</code>
A	<b>#Command</b>	<code>persp(rsm1, ~x1 + x2, zlab = "Rendimiento", col = rainbow(50), box = T, mgp = c(0.25, 0, 0), lty = 1, lwd = 1, font.axis = 2, font.lab = 2, border = "darkblue", cex.lab = 1.0, cex.axis = 0.9, contours = ("colors"))</code>
	<b>#Name</b>	<code>title(main = "A")</code>
B	<b>#Command</b>	<code>persp(rsm1, ~x1 + x3, zlab = "Rendimiento", col = rainbow(50), box = T, mgp = c(0.25, 0, 0), lty = 1, lwd = 2, font.axis = 2, font.lab = 2, border = "darkblue", cex.lab = 1.0, cex.axis = 0.9, contours = ("colors"))</code>
	<b>#Name</b>	<code>title(main = "B")</code>
C	<b>#Command</b>	<code>persp(rsm1, ~x1 + x4, zlab = "Rendimiento", col = rainbow(50), box = T, mgp = c(0.25, 0, 0), lty = 1, lwd = 2, font.axis = 2, font.lab = 2, border = "darkblue", cex.lab = 1.0, cex.axis = 0.9, contours = ("colors"))</code>
	<b>#Name</b>	<code>title(main = "C")</code>
D	<b>#Command</b>	<code>persp(rsm1, ~x2 + x3, zlab = "Rendimiento", col = rainbow(50), box = T, mgp = c(0.25, 0, 0), lty = 1, lwd = 2, font.axis = 2, font.lab = 2, border = "darkblue", cex.lab = 1.0, cex.axis = 0.9, contours = ("colors"))</code>
	<b>#Name</b>	<code>title(main = "D")</code>
E	<b>#Command</b>	<code>persp(rsm1, ~x2 + x4, zlab = "Rendimiento", col = rainbow(50), box = T, mgp = c(0.25, 0, 0), lty = 1, lwd = 2, font.axis = 2, font.lab = 2, border = "darkblue", cex.lab = 1.0, cex.axis = 0.9, contours = ("colors"))</code>
	<b>#Name</b>	<code>title(main = "E")</code>
F	<b>#Command</b>	<code>persp(rsm1, ~x3 + x4, zlab = "Rendimiento", col = rainbow(50), box = T, mgp = c(0.25, 0, 0), lty = 1, lwd = 2, font.axis = 2, font.lab = 2, border = "darkblue", cex.lab = 1.0, cex.axis = 0.9, contours = ("colors"))</code>
	<b>#Name</b>	<code>title(main = "F")</code>

**Figure 8.** Response Surface Graphics A and B

experimental region; moreover, the eigenvalues were negative (-0.676, -0.980, -1.729 and -2.261), indicating that the stationary point was at its maximum.<sup>31</sup> This is the kind of situation you aim for in response surface experimentation—clear evidence of a close set of optimal conditions.<sup>32</sup> For certification, a confirmatory experiment must be carried out close to this estimated optimum, in real values, at  $\text{pH} \approx 2.3$ , temperature  $\approx 74.7$  °C, time  $\approx 36.2$  and ultrasound power  $\approx 3.7$ .

Canonical analysis allows us to control the behavior of a second-order response surface, the understanding of

which is facilitated by maximum slope optimization through contour plots.<sup>25,27</sup> It is noteworthy that once it is concluded that in the experimentation region, the response variation is well modeled by a linear function of the control factors, then a search procedure for the best operating conditions can be started; that is, those levels of the quantitative control factors that optimize the response of interest.<sup>25,30,33</sup>

The algorithms used for maximum slope path optimization are described in more detail in the Supplementary material (Figure 39S and Figure 40S).

## 4. Conclusions

In this tutorial, the development of algorithms for data analysis in R language was presented. In addition, the main aspect of experimental planning in an unprecedented environment, R, was demonstrated through an original and simple experiment that involved the preparation of samples for the extraction of pectin from seriguela.

The performance of the proposed algorithm in R language was compared to that of other commercial and paid *software* packages that can be used in routine handling experiments, and the R language, in addition to being a free tool, was presented in this tutorial for easy handling.

Similar comparisons can be made using multivariate analysis. Thus, the authors understand that this tutorial, available in free *software*, will be a tool that will enable undergraduate and graduate students and researchers to develop data analysis and statistical forecasts easily and free of charge.

## Supplementary Data

Information on the R software used, R Studio, installation and loading packages, supplementary tables, and figures is provided in the supplemental material. Furthermore, the folder with the data used to reproduce the experiments carried out in this work can be downloaded through the following link:

Data

<https://drive.google.com/uc?export=download&id=1U94BkGNTVUhcPIBBTFBO192H8tFVhMNC>.

Supplementary Material

<https://drive.google.com/uc?export=download&id=1oTRMqhe207MxKPP7xEVAYWE62MADieCP>

## Acknowledgements

The authors are grateful to Coordenação de Aperfeiçoamento de Pessoal de Nível Superior, CAPES, Brazil (Code 001). We also thank CNPq (Grants 307429/2018-0 and 408767/2021-9). LBSF would like to thank the Federal Institute of Piauí (IFPI) for their encouragement and support for his doctorate. Additionally, LBSF thanks Dr. H.S. Barbosa and Prof. E.C. Muniz for their encouragement and guidance.

## Credit Authorship Contribution Statement

Luiz B. S. Filho: Conceptualization, Methodology, Investigation, Formal analysis, Data curatorship and

Writing – original draft; Ronaldo C. Coelho: Methodology, Formal analysis; Tiago L. S. Coêlho: Methodology, Formal analysis; Edvani C. Muniz: Validation, Writing – Review and Editing; Herbert de S. Barbosa: Supervision, Writing – Review and Editing.

## Declarations

The authors declare that there are no conflicts of interest of any type in this work or that the financial support for this work has influenced our findings.

## Bibliographic References

- Neto, B. B.; Ieda, S. S.; Bruns, R. E.; 25 anos de quimiometria no Brasil. *Quimica Nova* **2006**, *29*, 1401. [[Link](#)]
- Weissman, S. A.; Anderson, N. G.; Design of Experiments (DoE) and Process Optimization. A Review of Recent Publications. *Organic Process Research & Development* **2014**, *19*, 1605. [[Crossref](#)]
- Mäkelä, M.; Experimental design and response surface methodology in energy applications: A tutorial review. *Energy Conversion Management* **2017**, *151*, 630. [[Crossref](#)]
- Ebrahimi-Najafabadi, H.; Leardi, R.; Jalali-Heravi, M.; Experimental Design in Analytical Chemistry—Part I: Theory. *Journal of AOAC INTERNATIONAL* **2014**, *97*, 3. [[Crossref](#)]
- Garud, S. S.; Karimi, I. A.; Kraft, M.; Design of computer experiments: A review. *Computers & Chemical Engineering* **2017**, *106*, 71. [[Crossref](#)]
- Maran, J. P.; Priya, B.; Al-Dhabi, N. A.; Ponmurugan, K.; Moorthy, I. G.; Sivarajasekar, N.; Ultrasound assisted citric acid mediated pectin extraction from industrial waste of Musa balbisiana. *Ultrasonics Sonochemistry* **2017**, *35*, 204. [[Crossref](#)]
- Mugwagwa, L. R.; Chimphango, A. F. A.; Box-Behnken design based multi-objective optimization of sequential extraction of pectin and anthocyanins from mango peels. *Carbohydrate Polymers* **2019**, *219*, 29. [[Crossref](#)]
- Pinkowska, H.; Krzywonos, M.; Wolak, P.; Złocinska, A.; Pectin and Neutral Monosaccharides Production during the Simultaneous Hydrothermal Extraction of Waste Biomass from Refining of Sugar—Optimization with the Use of Doehlert Design. *Molecules* **2019**, *24*. [[Crossref](#)]
- Dejaegher, B.; Vander Heyden, Y.; Experimental designs and their recent advances in set-up, data interpretation, and analytical applications. *Journal of Pharmaceutical and Biomedical Analysis* **2011**, *56*, 141. [[Crossref](#)]
- Narendaran, S. T.; Meyyanathan, S. N.; Karri, V. V. S. R.; Experimental design in pesticide extraction methods: A review. *Food Chemistry* **2019**, *289*, 384. [[Crossref](#)]
- Coelho, T. L. S.; Braga, F. M. S.; Silva, N. M. C.; Dantas, C.; Lopes Júnior, C. A.; de Sousa, S. A. A.; Vieira, E. C.; Optimization of the protein extraction method of goat meat using factorial design and response surface methodology. *Food Chemistry* **2019**, *281*, 63. [[Crossref](#)]

12. Teófilo, R. F.; Ferreira, M. M.; Quimiometria II: planilhas eletrônicas para cálculos de planejamentos experimentais, um tutorial. *Química Nova* **2006**, *29*, 338. [[Link](#)]
13. Vriesmann, L. C.; Teófilo, R. F.; Petkowicz, C. L. D. O.; Optimization of nitric acid-mediated extraction of pectin from cacao pod husks (*Theobroma cacao* L.) using response surface methodology. *Carbohydrate Polymers* **2011**, *84*, 1230. [[Crossref](#)]
14. Minjares-Fuentes, R.; Femenia, A.; Garau, M. C.; Meza-Velázquez, J. A.; Simal, S.; Rosselló, C.; Ultrasound-assisted extraction of pectins from grape pomace using citric acid: A response surface methodology approach. *Carbohydrate Polymers* **2014**, *106*, 179. [[Crossref](#)]
15. Ezzati, S.; Ayaseh, A.; Ghanbarzadeh, B.; Heshmati, M. K.; Pectin from sunflower by-product: Optimization of ultrasound-assisted extraction, characterization, and functional analysis. *International Journal of Biological Macromolecules* **2020**, *165*, 776. [[Crossref](#)]
16. Shivamathi, C. S.; Gunaseelan, S.; Soosai, M. R.; Vignesh, N. S.; Varalakshmi, P.; Kumar, R. S.; Karthikumar, S.; Kumar, R. V.; Baskar, R.; Rigby, S. P.; Syed, A.; Elgorban, A. M.; Ganesh Moorthy, I. M.; Process optimization and characterization of pectin derived from underexploited pineapple peel biowaste as a value-added product *Food Hydrocolloids* **2022**, *123*. [[Crossref](#)]
17. Colodel, C.; Vriesmann, L. C.; Teófilo, R. F.; de Oliveira Petkowicz, C. L.; Extraction of pectin from ponkan (*Citrus reticulata* Blanco cv. Ponkan) peel: Optimization and structural characterization *International Journal of Biological Macromolecules* **2018**, *117*, 385. [[Crossref](#)]
18. Yang, J. S.; Mu, T. H.; Ma, M. M.; Optimization of ultrasound-microwave assisted acid extraction of pectin from potato pulp by response surface methodology and its characterization *Food Chemistry* **2019**, *289*, 351. [[Crossref](#)]
19. Breitreitz, M. C.; De Souza, A. M.; Poppi, R. J.; A didactic chemometrics experiment for design of experiments (DOE): evaluation of experimental conditions in the spectrophotometric determination of Iron II witho-phenanthroline. A tutorial, part III. *Química Nova* **2014**, *37*, 564. [[Crossref](#)]
20. Hilário, F. F.; Castro, J. P.; Barros, T. E.; Pereira-Filho, E. R.; Planejamento de misturas e visualização da região ótima com planilhas no excel: um Tutorial. *Química Nova* **2021**, *44*, 874. [[Crossref](#)]
21. Pereira, F. M. V.; Pereira-Filho, E. R.; Aplicação de programa computacional livre em planejamento de experimentos: um Tutorial. *Química Nova* **2018**, *41*, 1061. [[Crossref](#)]
22. Antony, J.; *Design of Experiments for Engineers and Scientists*, 2a. ed.; Elsevier, **2014**. [[Link](#)]
23. Ferreira, S. L. C.; Bruns, R. E.; Ferreira, H. S.; Matos, G. D.; David, J. M.; Brandão, G. C.; da Silva, E. G. P.; Portugal, L. A.; dos Reis, P. S.; Souza, A. S.; dos Santos, W. N. L.; Box-Behnken design: An alternative for the optimization of analytical methods. *Analytica Chimica Acta* **2007**, *597*, 179. [[Crossref](#)]
24. Neto, B. B.; Scarminio, I. S.; Bruns, R. E.; *Como fazer experimentos*, 4th. ed.; Bookman, **2010**. [[Link](#)]
25. Montgomery, D. C.; *Design and Analysis of Experiments*, 8th. ed.; John Wiley & Sons, **2012**. [[Link](#)]
26. Box, G. E. P.; Draper, N. R.; *Empirical model-building and response surfaces* New York, **1987**. [[Link](#)]
27. Myers, R. H.; Montgomery, D. C.; Anderson-Cook, C. M.; *Response Surface Methodology*, 3th. ed.; Wiley: Danvers, **2009**. [[Link](#)]
28. Guandalini, B. B. V.; Rodrigues, N. P.; Marczak, L. D. F.; Sequential extraction of phenolics and pectin from mango peel assisted by ultrasound. *Food Research International* **2019**, *119*, 455. [[Crossref](#)]
29. Sengar, A. S.; Rawson, A.; Muthiah, M.; Kalakandan, S. K.; Comparison of different ultrasound assisted extraction techniques for pectin from tomato processing waste. *Ultrasonics Sonochemistry* **2020**, *61*, 104812. [[Crossref](#)]
30. Lawson, J.; *Design and Analysis of Experiments with R* CRC Press - Taylor & Francis Group: Utah, **2015**. [[Link](#)]
31. Calado, V.; Montgomery, D.; *Planejamento de Experimentos usando o Statistica* E-Papers Serviços Editoriais LTDA: Rio de Janeiro, **2003**. [[Link](#)]
32. Bas, D.; Boyaci, I. H.; Modeling and optimization I: Usability of response surface methodology. *Journal of Food Engineering* **2007**, *78*, 836. [[Crossref](#)]
33. Lee, D. H.; Kim, S. H.; Byun, J. H.; A method of steepest ascent for multiresponse surface optimization using a desirability function method. *Quality Reliability Engineering International* **2020**, *36*, 1931. [[Crossref](#)]