#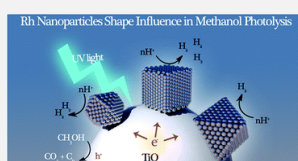 The TEM images-based predictive modeling for differently shaped Rh nanoparticles classification in a hybrid photocatalyst

Rh Nanoparticles Shape Influence in Methanol Photolysis

UFRGS — UNIVERSIDADE FEDERAL DO RIO GRANDE DO SUL

**M.A. Nazarkovsky[1], B. Albuquerque[2], G. Chacón[2], J. Dupont[2]**

[1]Chemistry Department, Pontifical Catholic University of Rio de Janeiro, 225 Marques de Sao Vicente Str., Rio de Janeiro 22451-900, Brazil
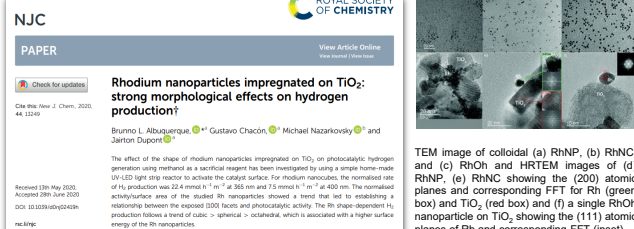[2]Laboratory of Molecular Catalysis, Institute of Chemistry, UFRGS, 9500 Bento Gonçalves Av., Porto Alegre 91501-970, RS, Brazil

**#DATA SCIENCE  #INDUSTRY 4.0  #PREDICTIVE MODELING  #MACHINE LEARNING  #AI  for  #CHEMISTRY**

## Development of the four materials science and engineering paradigms

- 1st – Empirical Science (experiments)
- 2nd – Theoretical Science (laws of natural sciences)
- 3rd – Computational Science, Simulations (DFT, molecular dynamics)
- **4th – Big Data-Driven Science (Artificial Intelligence)**

### IN THE BEGINNING....

ROYAL SOCIETY OF CHEMISTRY

NJC

PAPER

View Article Online

**Rhodium nanoparticles impregnated on TiO₂: strong morphological effects on hydrogen production†**

Bruno L. Albuquerque, *Gustavo Chacón, *Michael Nazarkovsky *and Jairton Dupont

Cite this: New J. Chem., 2020, 44, 13249

Received 13th May 2020,
Accepted 28th June 2020

DOI: 10.1039/d0nj02419h

rsc.li/njc

The effect of the shape of rhodium nanoparticles impregnated on TiO₂ on photocatalytic hydrogen generation using methanol as a sacrificial reagent has been investigated by using a simple home-made UV-LED light strip reactor to activate the catalyst surface. For rhodium nanocubes, the normalised rate of H₂ production was 22.4 mmol h⁻¹ m⁻² at 365 nm and 7.5 mmol h⁻¹ m⁻² at 400 nm. The normalised activity/surface area of the studied Rh nanoparticles showed a trend that led to establishing a relationship between the exposed (100) facets and photocatalytic activity. The Rh shape-dependent H₂ production follows a trend of cubic > spherical > octahedral, which is associated with a higher surface energy of the Rh nanoparticles.

TEM image of colloidal (a) RhNP, (b) RhNC, and (c) RhOh and HRTEM images of (d) RhNP, (e) RhNC showing the (200) atomic planes and corresponding FFT for Rh (green box) and TiO₂ (red box) and (f) a single RhOh nanoparticle on TiO₂ showing the (111) atomic planes of Rh and corresponding FFT (inset).

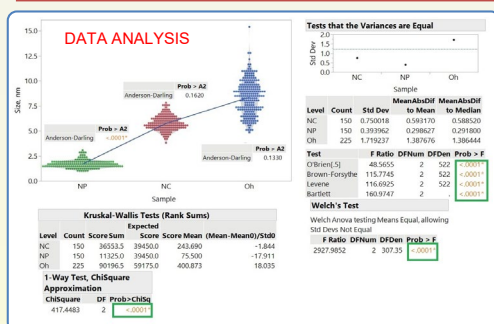### ...AND THE "WOW" IDEA EMERGED...

The results allow us to distinguish each shape by size distribution profiles from respective TEM microphotographs and make predictive modeling by means of machine learning algorithms!

### METHODOLOGY
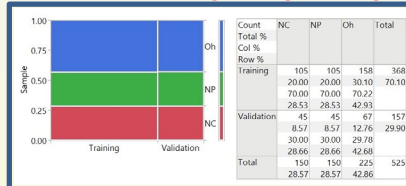
1) data analysis (the size distribution profiles distributions classification, analysis of variances, discriminant analysis);
2) machine learning with training (70%) and validation (30%) of the models stratified by the samples' types;
3) coded in JSL (JMP scripting language) associable with R, Python, Matlab and SAS;
4) scripts, calculators and some other details are available on the repository: https://github.com/Nazarkovsky/Rh-TiO2-classificator
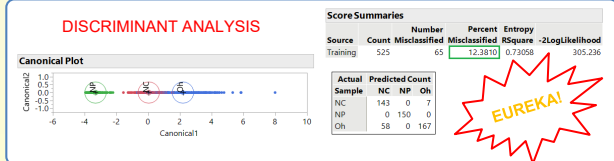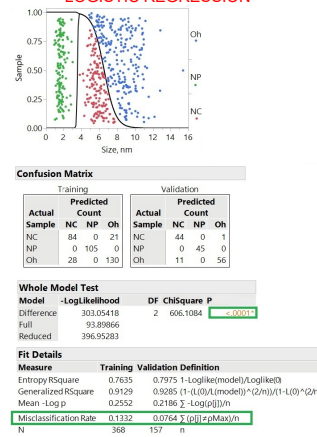
### THE MODELS:
K-Nearest Neighbors (KNN) – K = 100, Euclidean distances between the points
Bootstrap Forest - 1 split per sample, learning rate 0.1, 35 trees
Logistic Regression
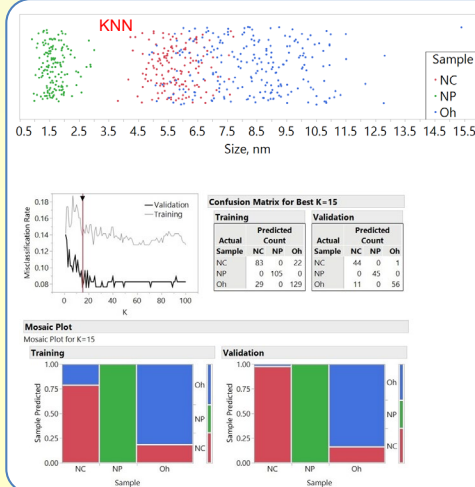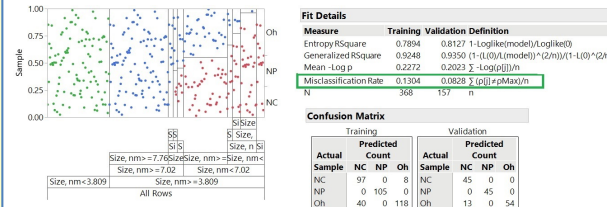Classification Tree - 9 splits
Naïve Bayes Classificator



## DATA ANALYSIS



## DISCRIMINANT ANALYSIS

EUREKA!



## TRAIN-VALIDATION BY SAMPLES
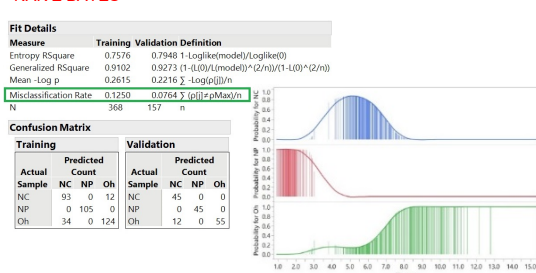


## LOGISTIC REGRESSION



## BOOTSTRAP FOREST



## KNN



## CLASSIFICATION TREE



## NAIVE BAYES

**Conclusions.** NC and Oh are normally distributed by the Anderson-Darling criterion, except NP whose nanoparticles size range is the narrowest among three samples. The variances are revealed to be heteroscedastic by all four tests, the non-parametric Kruskall-Wallis test has shown the non-equality for all three samples. The discriminant analysis at the overall MR of 12.38% has become promising to develop the machine learning algorithms for practical digital recognition of the samples by size. The most precise model is Logistic Regression at the misclassification rate MR is 7.64% with other better metrics (Generalized R² and Entropy R²) than another two models have at the same MR (Naïve Bayes and Bootstrap Forest). As a result, an offline HTML-calculator for Logistic Regression was developed.