# REGGAE

Regression Generator & Analyzer

How-to-reggae at Coding Camp
By Liliana Gallegos

# What is REGGAE?



Compiled data
(csv format)

Select analysis
Run regression

✓ Analysis report
✓ Regression plot

Aims:
 *Statistical reproducibility*
 *open-source*
 *user-friendly*
 *quick analysis*

# Setup

❏ Create r environment and install r-essentials:

    `conda create –n r_env r-essentials r-base`

❏ Activate environment:

    `conda activate r_env`

❏ Install packages required:

    `conda install --yes --file R-requirements.txt`

❏ Confirm by running reggae help options:

    `Rscript reggae.r -h`

**Available:** *https://github.com/Liliana-Gallegos/REGGAE*

# Statistical analysis options

Regression analysis: Multivariate Linear or Random Forest (*quick not opt*)

| | |
|---|---|
| –m , --model | o Feature selection: *full, stepwise, dredge, mincorr* |
| –b , --buildmodel | o Build a model: $x_1, x_2, x_n$ |
| –y , --yresponse | o Define y-response |
| | o Split into Train/ Test datasets: [Default = 1 full dataset] |
| –r , --randsample | o (a) ratio between 0 and 1 or (b) 0 = predefined test/train |
| –p , --pca | o Principal Component analysis with kmeans clustering. |
| | o Build scaled/unscaled train and test sets random or universal training set. |
| –c , --corrplot | o Pairwise Correlation |
| –q , --crossvalidation | o Cross Validation $q^2$ values from leave-one-out, K-fold, and external. |
| –d , --diagnostics | o Diagnostics QSAR criteria for an acceptable model, F-value comparisons, collinearity, *outlier testing*. |
| | o ANOVA analysis. |
| –v , --verbose | o Plots for selected analysis. |

# Other options

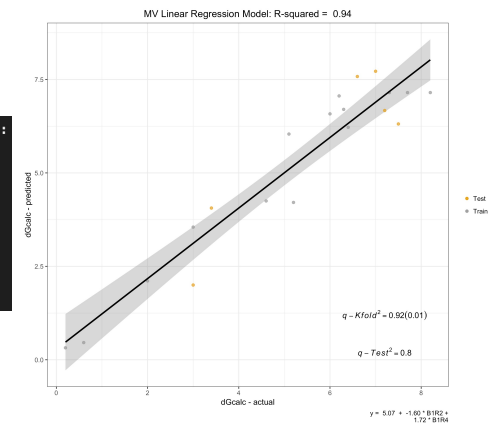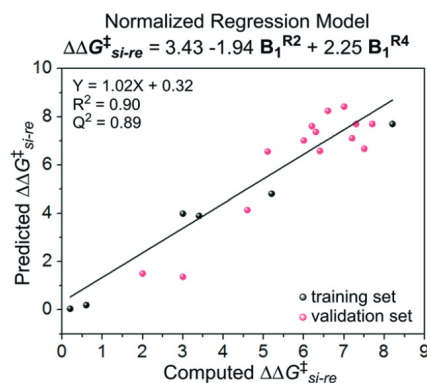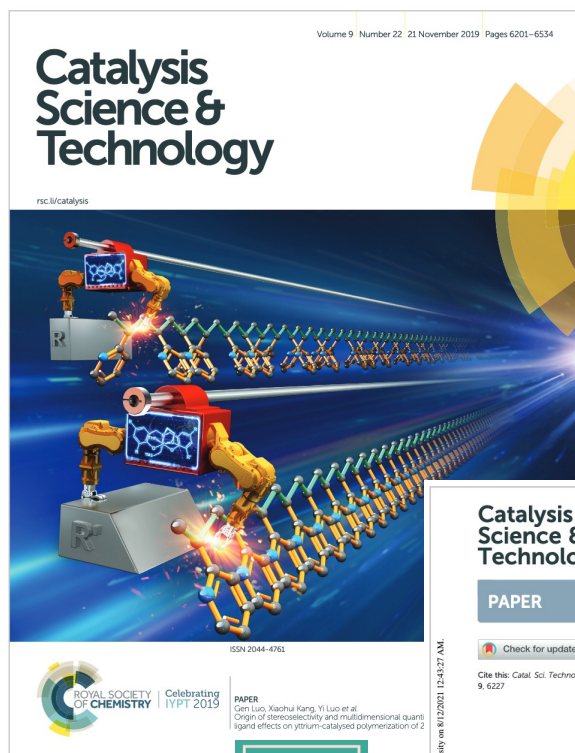| | |
|---|---|
| –i , --inputfile | Requires csv file name. |
| –o , --outputfile | Optional: output file name (default = REGGAE-analysis-output.txt) |
| –x , --extdata | To validate model using external csv data set. Requires the name of csv data file. |
| –e , --exportdata | Export data sets. Select from option: scaled, unscaled, predicted. |
| –s , --seed | Optional: specify the seed for random sample split. (default seed = 42) |
| –k , --kfoldvalue | Optional: To adjust k-value for K-fold cross validation. (default = 5) |
| –v , --verbose | Optional: To print extra output and plots. (defaul = FALSE) |
| –h , --help | |

Rscript reggae.r –i data.csv –m stepwise –y dG –r 0.8 –q –d –v

Predictive Model & Report

MV Linear Regression Model: R-squared = 0.94

o Linear regression model with SELECTED features:
  Number of features (including response):  3

a) scaled coefficients:
   y =  5.073  +  -1.602 * B1R2 + 1.720 * B1R4

b) unscaled:
   y = 4.23 + -2.92 * B1R2 + 1.61 * B1R4

  R2-train =  0.94
  adj R2-train =  0.93
   RMSE-train =  0.58

# Example: Quantitative Structure–Selectivity Relationships



### Normalized Regression Model

$$\Delta\Delta G^{\ddagger}_{si\text{-}re} = 3.43 - 1.94\, B_1^{R2} + 2.25\, B_1^{R4}$$

$Y = 1.02X + 0.32$
$R^2 = 0.90$
$Q^2 = 0.89$

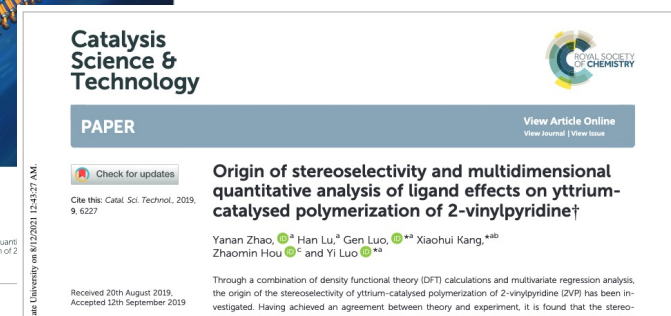- training set
- validation set

Total of 21 catalysts
Steric and electronic features:
○ Sterimol descriptors:
   $B_1$, $B_5$, $L$ and %Vbur
○ Natural Population Analysis (NPA) charges

For statistical reproducibility:
○ Available dataset (ESI: *.xlsx, .csv, table*)
○ Available code
○ Labeled split datasets or methods used



| | $R^2$ | | | $R^4$ | | | NPA Charge | | | %V_Bur | $\Delta G_{DFT}$ | $\Delta G_{predicted}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $L^{R2}$ | $B_1^{R2}$ | $B_5^{R2}$ | $L^{R4}$ | $B_1^{R4}$ | $B_5^{R4}$ | $Q_Y$ | $Q_{O1}$ | $Q_{O2}$ | | | |
| A | 4.62 | 2.98 | 6.07 | 4.62 | 2.98 | 6.07 | 1.92 | -0.92 | -0.92 | 89.7 | 0.2 | 0.0 |
| B | 4.41 | 2.93 | 3.35 | 4.62 | 2.98 | 6.07 | 1.91 | -0.92 | -0.92 | 89.7 | 0.6 | 0.2 |
| C | 3.09 | 1.70 | 2.19 | 4.62 | 2.98 | 6.07 | 1.93 | -0.91 | -0.90 | 88.7 | 3.4 | 3.9 |

```
Rscript reggae.r –i 2vp-labeled.csv –b B1R2,B1R4 –y dGcalc –r 0 –q –d –v
```

Zhao, Y.; Lu, H.; Luo, G.; Kang, X.; Hou, Z.; Luo, Y. *Catal. Sci. Technol.,* **2019**, *9,* 6227–6233.
Krafczyk, M. S.; Shi, A.; Bhaskar, A.; Marinov, D. Stodden, V. *Phil. Trans. R. Soc. A.* **2021,** *379.*