3rd Internal CINE-CMSC Workshop

# Data Mining Tools Applied to Quantum Chemistry Data

Speaker:
Johnatan Mucelini

Supervisor:
Prof. Dr. Juarez L. F. Da Silva

QC Collaborators:
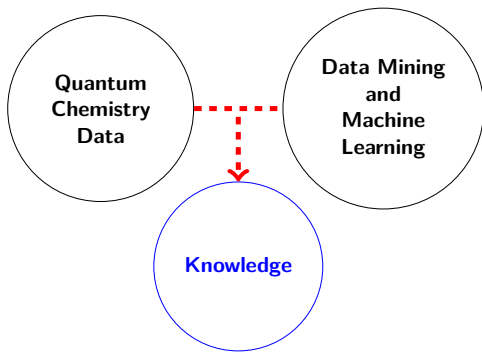Paulo de Carvalho Dias Mendes,
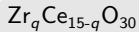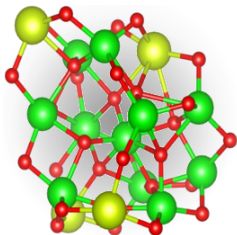Priscilla Felício-Sousa,
Karla F. Andriani.
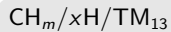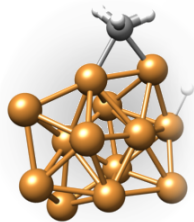
DM/ML Collaborators:
Marcos G. Quiles,
Ronaldo C. Prati.

# Introduction - Materials



$Zr_q Ce_{15-q} O_{30}$

• Employed in TWC and candidate for others process.
• Explore material in nanoparticles structures.

$Pt_l TM_{55-l}$

• Pt-based Catalysts are wildly employed.
• Explore this alloys structural preferences.

$CH_m/xH/TM_{13}$

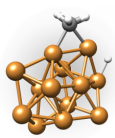• $CH_4$ dehydrogenation candidates.
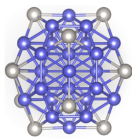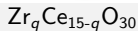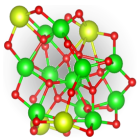• Reaction intermediates study.

Motivation:

- Develop the chem/comp overlap;
- Develop material science area.

Objectives:

- Find new patterns;
- Initial Studies;
- Tools Development.

# Data Mining - Challenge

**Calculations Sets:**



| | $Zr_qCe_{15-q}O_{30}$ | $Pt_lTM_{55-l}$ | $CH_m/xH/TM_{13}$ |
|---|---|---|---|
| $n$ | 1600 | 330 | 700 |
| vars | $q$ (16) | $l$ (2) | $(m,x)$ (9) |
| | | TM (9) | TM (4) |

**Initial Data:**

Molecular info (structured data):
- Energy $[-42.0\,eV]$;
- HOMO $[-4.123\,eV]$;

Atomic info (attribute-vector):
- Exposed to vacuum [True, False, ...];
- Atomic Charges [0.7, -0.8, ...].

**DM/ML Input:**

| attribute-value table | feature0 | feature1 | ... |
|---|---|---|---|
| sample0 | data00 | data01 | ... |
| sample1 | data10 | data11 | ... |
| ... | ... | ... | ... |

# Data Mining - Challenge



**Calculations Sets:**

| | $Zr_qCe_{15-q}O_{30}$ | $Pt_lTM_{55-l}$ | $CH_m/xH/TM_{13}$ |
|---|---|---|---|
| $n$ | 1600 | 330 | 700 |
| vars | $q$ (16) | $l$ (2) | $(m, x)$ (9) |
| | | TM (9) | TM (4) |

Quandarium (python):

- Find Calc.          • Find surf.
- Extract Info.        atoms
- Geometical        • Connections
  Analysis:            Analysis
- $ECN/d_{av}$

It operate recursively!

**Initial Data:**

Molecular info (structured data):
- Energy $[-42.0\,eV]$;
- HOMO $[-4.123\,eV]$;

Atomic info (attribute-vector):
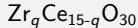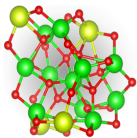- Exposed to vacuum [True, False, ...];
- Atomic Charges [0.7, -0.8, ...].

**DM/ML Input:**

| attribute-value table | feature0 | feature1 | ... |
|---|---|---|---|
| sample0 | data00 | data01 | ... |
| sample1 | data10 | data11 | ... |
| ... | ... | ... | ... |

# Data Mining - Challenge



**Calculations Sets:**

| | $Zr_qCe_{15-q}O_{30}$ | $Pt_lTM_{55-l}$ | $CH_m/xH/TM_{13}$ |
|---|---|---|---|
| $n$ | 1600 | 330 | 700 |
| vars | $q$ (16) | $l$ (2) | $(m, x)$ (9) |
| | | TM (9) | TM (4) |

Quandarium (python):

- Find Calc.          • Find surf.
- Extract Info.       atoms
- Geometical         • Connections
  Analysis:           Analysis
- $ECN/d_{av}$

It operate recursively!

**Initial Data:**

Molecular info (structured data):
- Energy $[-42.0\,eV]$;
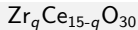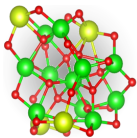- HOMO $[-4.123\,eV]$;

Atomic info (attribute-vector):
- Exposed to vacuum [True, False, ...];
- Atomic Charges [0.7, -0.8, ...].

**DM/ML Input:**

| attribute-value table | feature0 | feature1 | ... |
|---|---|---|---|
| sample0 | data00 | data01 | ... |
| sample1 | data10 | data11 | ... |
| ... | ... | ... | ... |

How to get molecular data from atomic data?

# Data Mining

## How to get molecular data from atomic data?

Take **molecular data** from **operator** over a **bag** (properties), for all or a **class** of atoms.

$$Molecular\_Data = \underbrace{OPERATOR}_{average} \left[ \underbrace{BAG}_{ECN} \left[ \underbrace{CLASS}_{Pt} \right] \right]$$

$$1.65 = Av. \left[ \begin{bmatrix} 1.1 \\ 2.2 \end{bmatrix} \right] = Av. \left[ \begin{bmatrix} 1.1 \\ 3.3 \\ 2.2 \\ 3.1 \end{bmatrix} \begin{bmatrix} ✔ \\ ✘ \\ ✔ \\ ✘ \end{bmatrix} \right]$$

**Operator**: Operates over one or more arguments, and return a number (Ex.: sum).

**Classes**: Set of atoms that meet a condition!
- Ex.: O, O exposed to the vacuum, O exposed to the vacuum with $1 < ECN < 2$, ....

# Data Mining

## How to get molecular data from atomic data?

Take **molecular data** from **operator** over a **bag** (properties), for all or a **class** of atoms.

$$Molecular\_Data = \underbrace{OPERATOR}_{\text{average}} \left[ \underbrace{BAG}_{\text{ECN}} \left[ \underbrace{CLASS}_{\text{Pt}} \right] \right]$$

$$1.65 = \text{Av.} \left[ \begin{bmatrix} 1.1 \\ 2.2 \end{bmatrix} \right] = \text{Av.} \left[ \begin{bmatrix} 1.1 \\ 3.3 \\ 2.2 \\ 3.1 \end{bmatrix} \begin{bmatrix} ✔ \\ ✘ \\ ✔ \\ ✘ \end{bmatrix} \right]$$

**Operator**: Operates over one or more arguments, and return a number (Ex.: sum).

**Classes**: Set of atoms that meet a condition!
- Ex.: O, O exposed to the vacuum, O exposed to the vacuum with $1 < ECN < 2$, ....

Quandarium (python):

Flexible data manipulation:
- bag → class
- class1 + class2 → class3
- class + bag → class
- class → bag

It operate recursively!

# Data Mining

## How to get molecular data from atomic data?

Take **molecular data** from **operator** over a **bag** (properties), for all or a **class** of atoms.

$$Molecular\_Data = \underbrace{OPERATOR}_{\text{average}} \left[ \underbrace{BAG}_{\text{ECN}} \left[ \underbrace{CLASS}_{\text{Pt}} \right] \right]$$

$$1.65 = \text{Av.} \left[ \begin{bmatrix} 1.1 \\ 2.2 \end{bmatrix} \right] = \text{Av.} \left[ \begin{bmatrix} 1.1 \\ 3.3 \\ 2.2 \\ 3.1 \end{bmatrix} \begin{bmatrix} \checkmark \\ \times \\ \checkmark \\ \times \end{bmatrix} \right]$$

**Operator**: Operates over one or more arguments, and return a number (Ex.: sum).

**Classes**: Set of atoms that meet a condition!
- Ex.: O, O exposed to the vacuum, O exposed to the vacuum with $1 < ECN < 2$, ....

Quandarium (python):

Flexible data manipulation:
- bag $\rightarrow$ class
- class1 + class2 $\rightarrow$ class3
- class + bag $\rightarrow$ class
- class $\rightarrow$ bag

It operate recursively!

Thus we access many molecular data (attribute-value table). Lets **analyse** the data!

# Correlation Analysis

### Pearson

$$r = \frac{cov(x, y)}{\sigma(x)\sigma(y)}$$

• Non-ranked data;
• Sensitive to outliers.

### Spearman

$$r_s = \frac{cov(r_x, r_y)}{\sigma(r_x)\sigma(r_y)}$$

• Ranked data;
• Robust to outliers.

Correlation Interpretation:

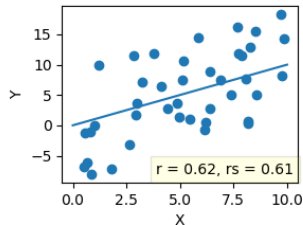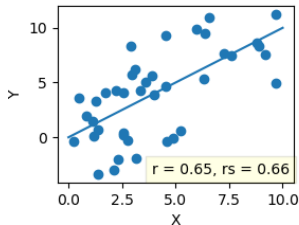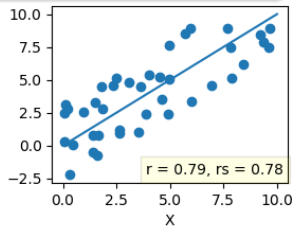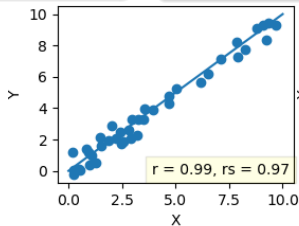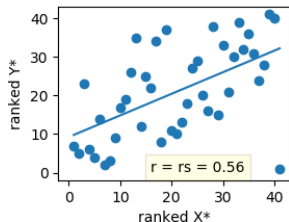- $-1 \geq r \geq 1$

- If $X$ increase, $Y$ is **expected** to go:

$$\overbrace{r < \underbrace{0}_{} < r}^{downward}$$

$\phantom{r < 0}$ upward

- How much expected?
  "How strong correlated?"

  "As large as was $|r|$."



r = 0.99, rs = 0.97

r = 0.79, rs = 0.78

r = 0.65, rs = 0.66

r = 0.62, rs = 0.61

# Correlation Analysis

### Pearson

$$r = \frac{cov(x, y)}{\sigma(x)\sigma(y)}$$

•Non-ranked data;
•Sensitive to outliers.

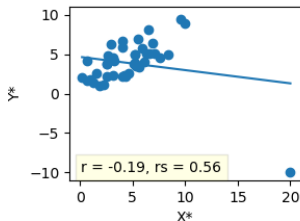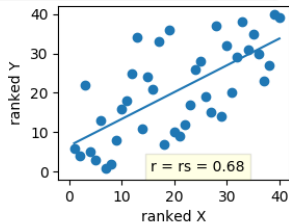### Spearman

$$r_s = \frac{cov(r_x, r_y)}{\sigma(r_x)\sigma(r_y)}$$

•Ranked data;
•Robust to outliers.



An example of outlier effect:
- $r$ reduced 0.89;
- $r_s$ reduced 0.12;

# Correlation Analysis

### Pearson

$$r = \frac{cov(x,y)}{\sigma(x)\sigma(y)}$$

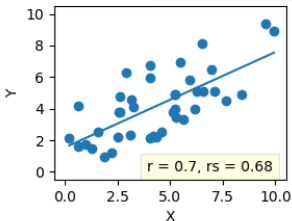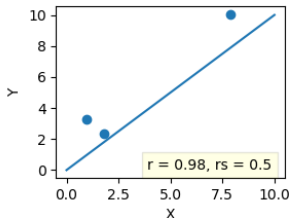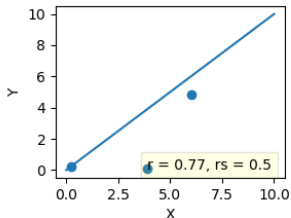- Non-ranked data;
- Sensitive to outliers.

### Spearman

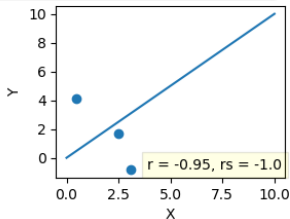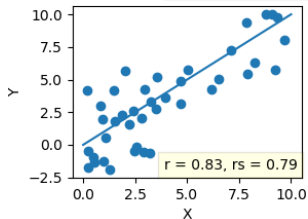$$r_s = \frac{cov(r_x, r_y)}{\sigma(r_x)\sigma(r_y)}$$
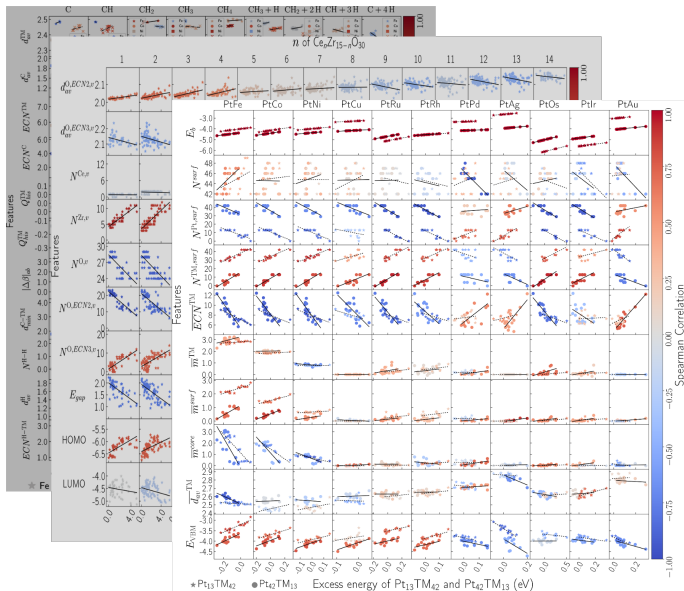
- Ranked data;
- Robust to outliers.

- Should I trust large correlation?

  Depend on data **size** and **distribution**.

  Good practice: Hypothesis test (Bootstrap) and pvalue.



r = 0.83, rs = 0.79

r = -0.95, rs = -1.0

r = 0.77, rs = 0.5

r = 0.98, rs = 0.5

# Data Representation



Scatter-plot Matrix:
- Rows: Features;
- Columns: Dataset part.

Cell:
- Scatter-plot;
- Y axis: Feature;
- X axis: Energy;
- Correlation: Colors;
- Linear Model.
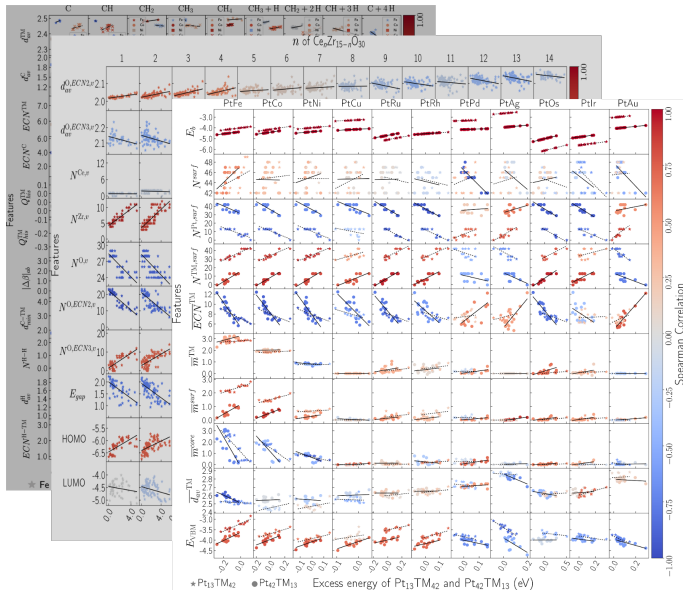
# Data Representation



Scatter-plot Matrix:
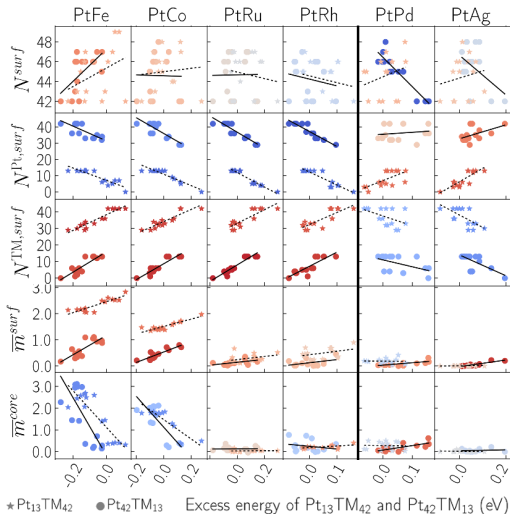- Rows: Features;
- Columns: Dataset part.

Cell:
- Scatter-plot;
- Y axis: Feature;
- X axis: Energy;
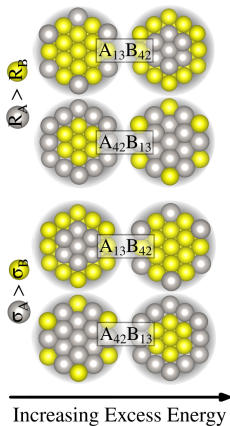- Correlation: Colors;
- Linear Model.

Quandarium:
- Scatternplot
- bag histograms

# Results

## Pt$_{13}$TM$_{42}$ and Pt$_{42}$TM$_{13}$



- High correlation! Few samples.
- TM *vs* Pt sites preference ($N$,$ECN$);
- Influence in other properties ($m$,$d_{av}$);

R: Radius = $d_{av}^{bulk}/2$

σ: Surface Energy

$x$: Electronegativity

R(Å)

| Fe | Co | Ni | Cu |
|------|------|------|------|
| 1.26 | 1.24 | 1.24 | 1.28 |
| Ru | Rh | Pd | Ag |
| 1.34 | 1.35 | 1.39 | 1.47 |
| Os | Ir | **Pt** | Au |
| 1.36 | 1.37 | **1.40** | 1.47 |

σ (eV/atom)

| Fe | Co | Ni | Cu |
|------|------|------|------|
| 0.88 | 0.71 | 0.65 | 0.47 |
| Ru | Rh | Pd | Ag |
| 1.05 | 0.81 | 0.56 | 0.33 |
| Os | Ir | **Pt** | Au |
| 1.21 | 0.90 | **0.64** | 0.32 |

Increasing Excess Energy

## Conclusion

## Quantum Chemistry Data Mining with Correlations:

Analysis Benefits:
- Easy access to useful **chemistry information**;
- **Quantitative trends** analysis;
- Very **little explored**.

Analysis Limitation:
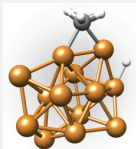- Small **dataset size**;
- Many **variables** in the study.

Analysis Applicability:
- Can be applied to **any material**;
- Require small **programming skills**;
- Require some **statistical concepts**.

# Perspectives 2019 - 2020

Complete the data mining works in progress!
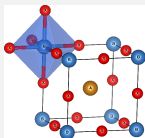
### Molecules over TM



- $CH_n/H_m/TM_{13}$ Dataset (previously presented).

### Perovskites



- Solid State Feature Extraction.

### Article: Nanocluster DM Analysis

- **DM** and **correlation analysis**;
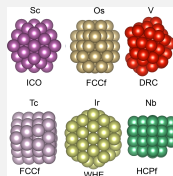- For **chemists**;
- 5 QC datasets (3+2);
- Release **Quandarium**.

TM Nanoclusters and Alloys Energy Regression:

### TM nanoclusters and alloys:



- Employ several previous QTNano studies ($TM_{13}$, $TM_{55}$);
- Methodological normalisation;
- Algorithms: MLP, random-forest, kernel regression...

# Acknowledgements



Thanks for your Attention!