# Topological Analysis and ML for biological docking

Gleb Rusyaev    Letovo

# Outline

# What is it?



Antibody

- Antibody vs Bad guys (Pathogens, Viruses, Cancer Cells, etc …)

- Biological (and physical) systems gravitate towards lowest potential energy

# Why bother?

1. Check if synthetic anti-cell can glue to the virus
2. Synthesize
3. Cure virus
4. ???
5. PROFIT

# How to calculate potential energy?

(i) Coulomb's law: $|F| = k_e \frac{|q_1 q_2|}{r^2}$

(ii) Van der Waals force

(iii) Energy of Statistic Potential

# Dataset review: docking_data.csv

(i) `E_fst`: Anti-Body Total Energy

(ii) `E_fst_elstat`: Anti-Body Coulomb's Energy

(iii) `E_fst_VdW`: Anti-Body Statistical Potential Energy

(iv) `group_size`: size of group (in which it included)

(v) `type_of_complex_in_group` :  center | in_between,
border ): place in group

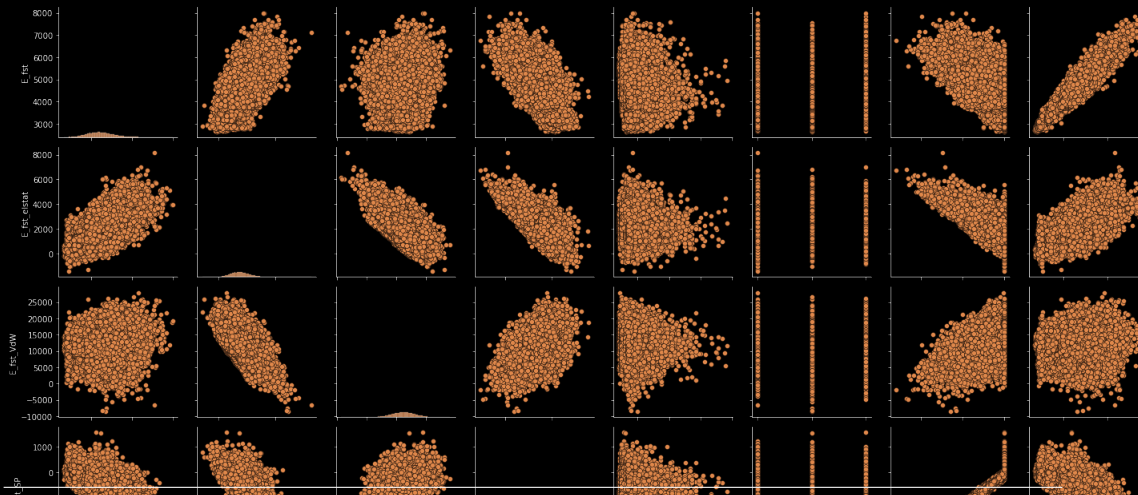(vi) `E_snd_SP`: SP of individual in group

# Dataset review: docking_data.csv

 (i) `avg_E_fst`: average energy in group
 (ii) `avg_E_snd_SP`: average stat. potential energy in group
(iii) `E_third`: Minimal Energy Complex after optimisation
(iv) `alt_E_third`: Minimal Energy Complex after optimisation (alternative method)
 (v) `E_third_SP`: Stat. Potential after optimisation
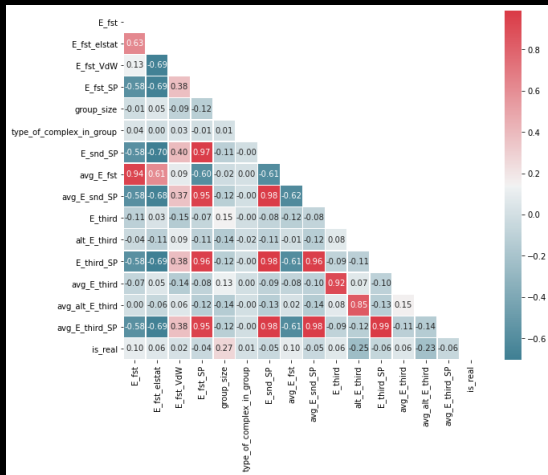(vi) `avg_E_third`: Avg. `E_third` in group

# Dataset review: docking_data.csv

  (i) `avg_alt_E_third`: Avg. `alt_E_third` in group
 (ii) `avg_E_third_SP`: Avg. `E_third_SP` in group
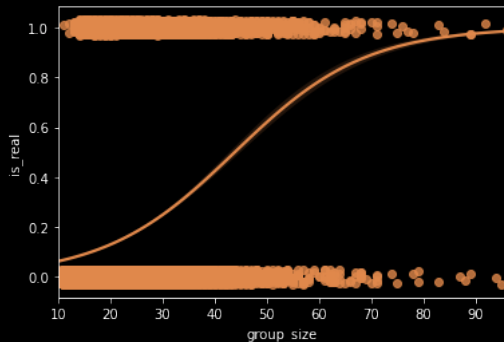(iii) `is_real`: Does it actually works this?
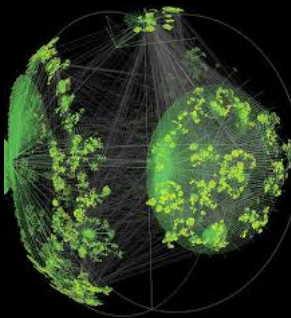
# Naive Correlation

# Naive Correlation

# Multi-Linear Regression



1. Complex border has 4x higher chance in being real, than centroid
2. Group size have 0.0127 correlation with `is_real` (as we early noticed, it's correlated)

# Topology



1337-dimensions $\xrightarrow{\text{hyper-cube clusterization}}$ nice 2d graph

# 10, 20, 40 Hyper-cubes

# 100 Hyper-cubes

Biological docking: 100 hypercubes

# Gradient Boosted Decision Trees

- Performance on training set:
  0.63 − 0.66
- Performance on validating set:
  0.6503
- Can be done even better using
  AWS

# References

Jupyter Notebook: `https://rusyaew.github.io/DockingML.ipynb`
Dataset: `https://rusyaew.github.io/docking_data.csv`
Bye!