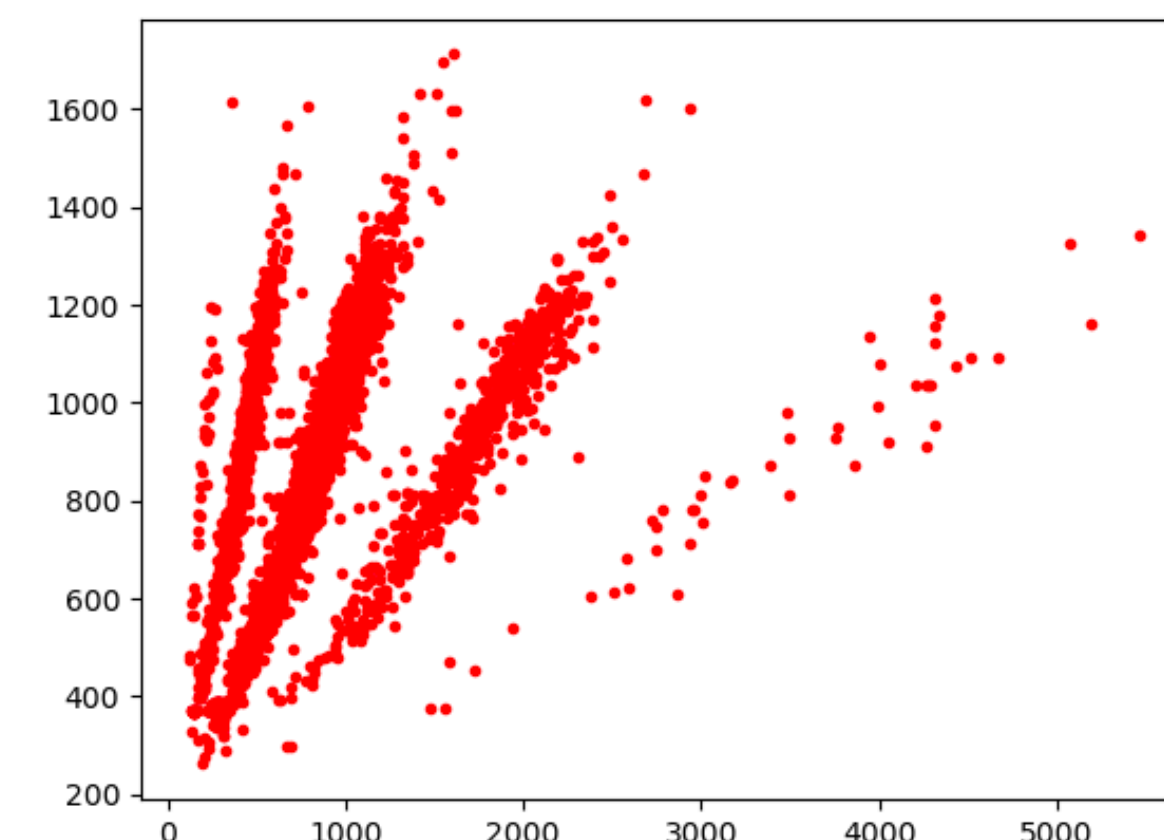# Condensed Matter Theory Group

**Tom, Gino, Milan, Peru, Chris**

**24 March 2021**

# Used SKlearn standard methods.



```python
40  # %% Train / test splitting
41  target = data['is_centrosymmetric']
42
43  X_train, X_test, y_train, y_test = train_test_split(
44      features, target, test_size=0.33, random_state=42)
45
46  y_train = y_train.to_numpy()
47
48  # %% Full model defn as pipeline
49  pclf = Pipeline([
50      ('imputer', SimpleImputer(strategy='mean', verbose=1)),
51      ('scaler', MinMaxScaler()),
52      ('feature_sel', SelectKBest(chi2, k = 50)),
53      ('fitting', RandomForestClassifier(random_state=0))
54  ])
55  # %% Fitting
56  pclf.fit(X_train, y_train)
57
58  # %% Prediction
59  y_pred = pclf.predict(X_test)
60  print('f1 score: ', f1_score(y_test, y_pred, average = 'macro'))
61  # %% testing
62
63  test_csvs = glob.glob("./data/test_*.csv")
64  tests = {Path(t).stem : pd.read_csv(t) for t in test_csvs}
65
66  test_data = pd.concat([
67      read_descriptors('./data/test_descriptors.csv'),
68      #tests['test_rdk'].drop('0', axis = 1),
69      tests['test_mord3d'].drop(['identifiers', 'Unnamed: 0', 'name', 'InchiKey',
      'smiles'], axis = 1),
70      #tests['test_mol2vec'],
71      ], axis = 1)
72
73  pclf.fit(features, target)
74  test_pred = pclf.predict(test_data)
```

- Random Forests for categorisation.

- Linear Kernal Ridge Regression on the regression.

- Used default sklearn feature scaling and imputers for cleanup.

- Refit on the full train data before doing a final prediction.

- Above: Z' in the cell volume.

# Other ideas

- Mol2vec: Unsupervised Machine Learning Approach with Chemical Intuition

- SchNet: A continuous-filter convolutional neural network for modeling quantum interactions