



UNIVERSITÀ DEGLI STUDI DI NAPOLI FEDERICO II - DIPARTIMENTO DI

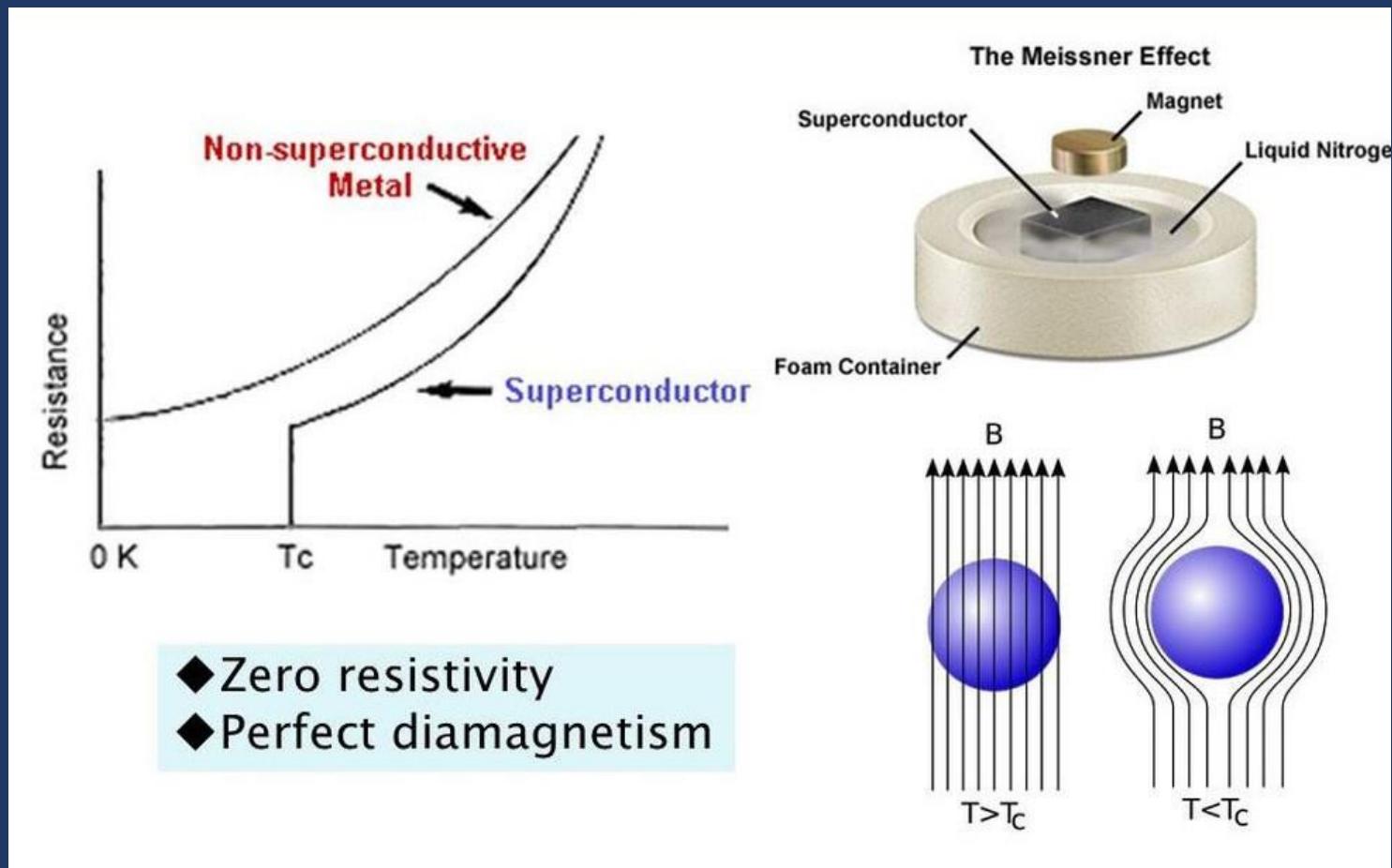
FISICA "ETTORE PANCINI"

# Machine Learning Approach for Prediction of Critical Temperature of Superconductor Materials



# What is Superconductivity?

- Discovered by Onnes in 1911
- Zero Resistivity, Meissner Effect, Fase Transition
- Low Temperature VS High Temperature Superconductors
- BCS Theory
- Applications: Quantum Computers, Magnetic Levitation, Electromagnets for Engineering and many others.



# Dataset Introduction

- **21263 Samples**
  - SuperCon Online Database of NIMS (Japan's National Institute for Materials Science)
  - Link: <https://archive.ics.uci.edu/ml/datasets/Superconductivity+Data>
  - We are going to study the data saved in the file "train.csv, which are the result calculations on the data of the file 'unique.csv'.
- **Continous Target Variable**
- **81 Feature Real Variables**
  1. Number of Elements
  2. Atomic Mass
  3. Atomic Radius
  4. First Ionization Energy
  5. Density
  6. Electron Affinity
  7. Fusion Heat
  8. Thermal Conductivity
  9. Valence

# Content

- **Exploratory Data Analysis**

- Missing and Categorical Data
- Scatter Plot Matrix and Correlation Matrix

- **Data Preprocessing**

- Feature Scaling
- Dimensionality Reduction
- Cross Validation

- **Evaluation and Tuning**

- Scoring
- Hyperparameter Tuning

- **Regression models**

- Linear Regression
- Non-linear Regression
- Ensemble

- **Conclusions**

- Results
- Running Time

# Exploratory Data Analysis

- No Missing Data
- No Categorical Data

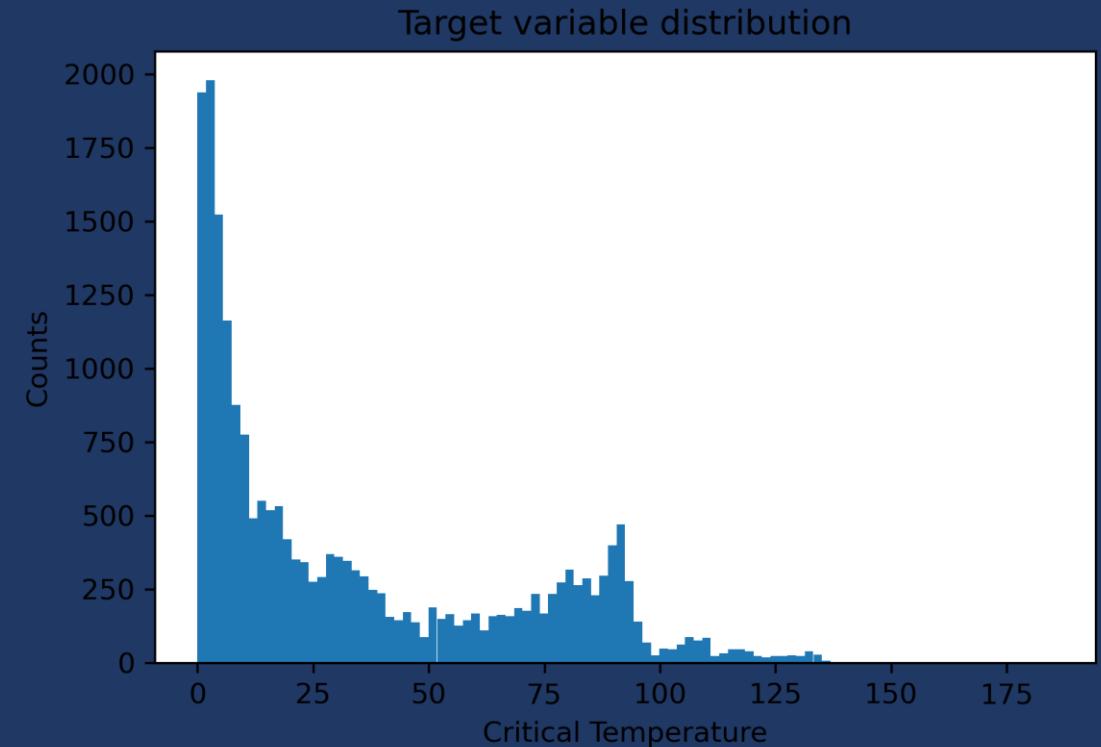
Stats	Formula
MEAN	$\mu = \sum_{i=1}^n \frac{t_i}{i}$
WEIGHTED MEAN	$\nu = \sqrt{\sum_{i=1}^n p_i t_i}$
GEOMETRIC MEAN	$= \prod_i^n t_i^{\frac{1}{n}}$
WEIGHTED GEOMETRIC MEAN	$= \prod_i^n t_i^{p_i}$

Stats	Formula
ENTROPY	$= - \sum_{i=1}^n w_i \log(w_i)$
WEIGHTED ENTROPY	$= - \sum_{i=1}^n A_i \log(A_i)$
RANGE	$t_{max} - t_{min}$
WEIGHTED RANGE	$p(t_{max})t_{max} - p(t_{min})t_{min}$
STANDARD DEVIATION	$\sqrt{\frac{1}{2} \sum_i^n (t_i - \mu)^2}$
WEIGHTED STANDARD DEVIATION	$\sqrt{p_i \sum_i^n (t_i - \mu)^2}$

# Exploratory Data Analysis

- Continuous Target Variable: Critical Temperature
  - Regression problem

STATS	CRITICAL TEMPERATURE
MEAN	34,4
STD	34,2
MIN	0,0
MAX	185,0



# Correlation Matrices

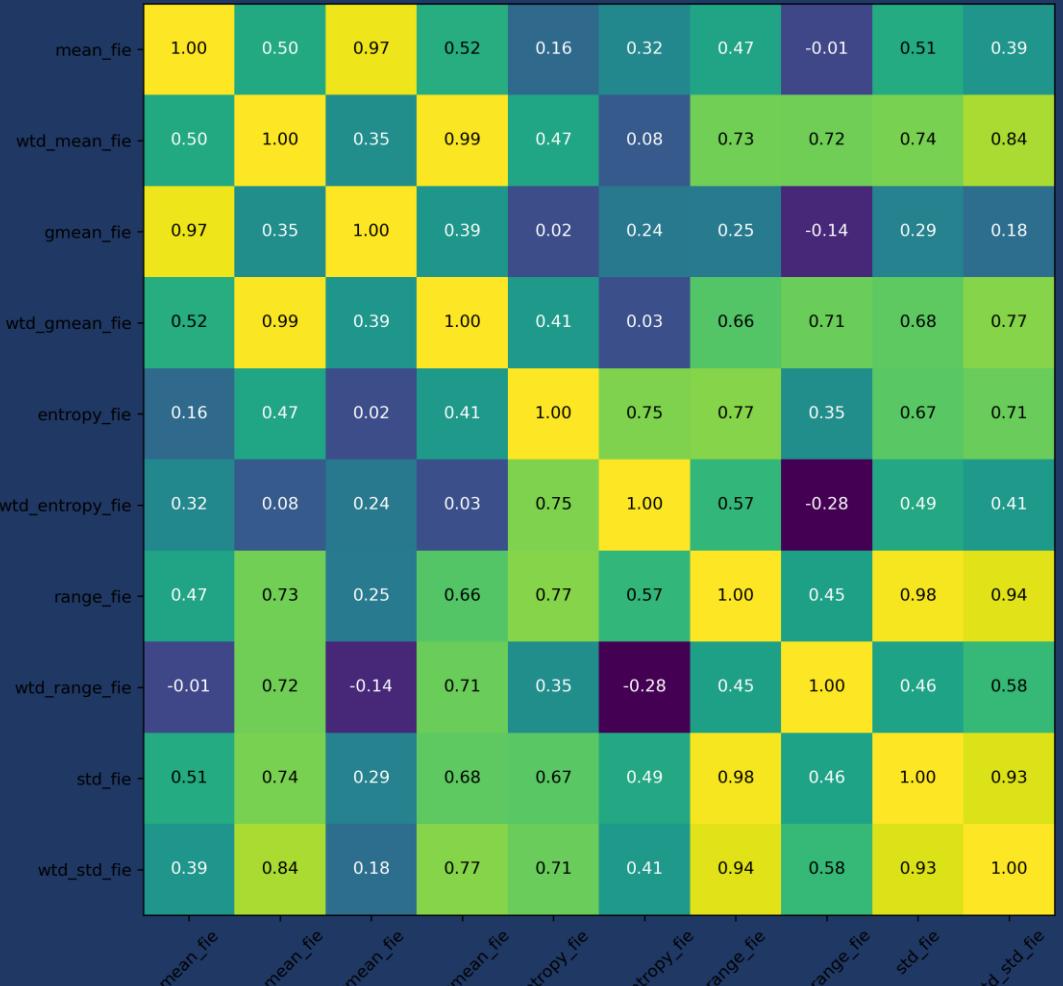
Pearson Correlation Coefficients:

$$r = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

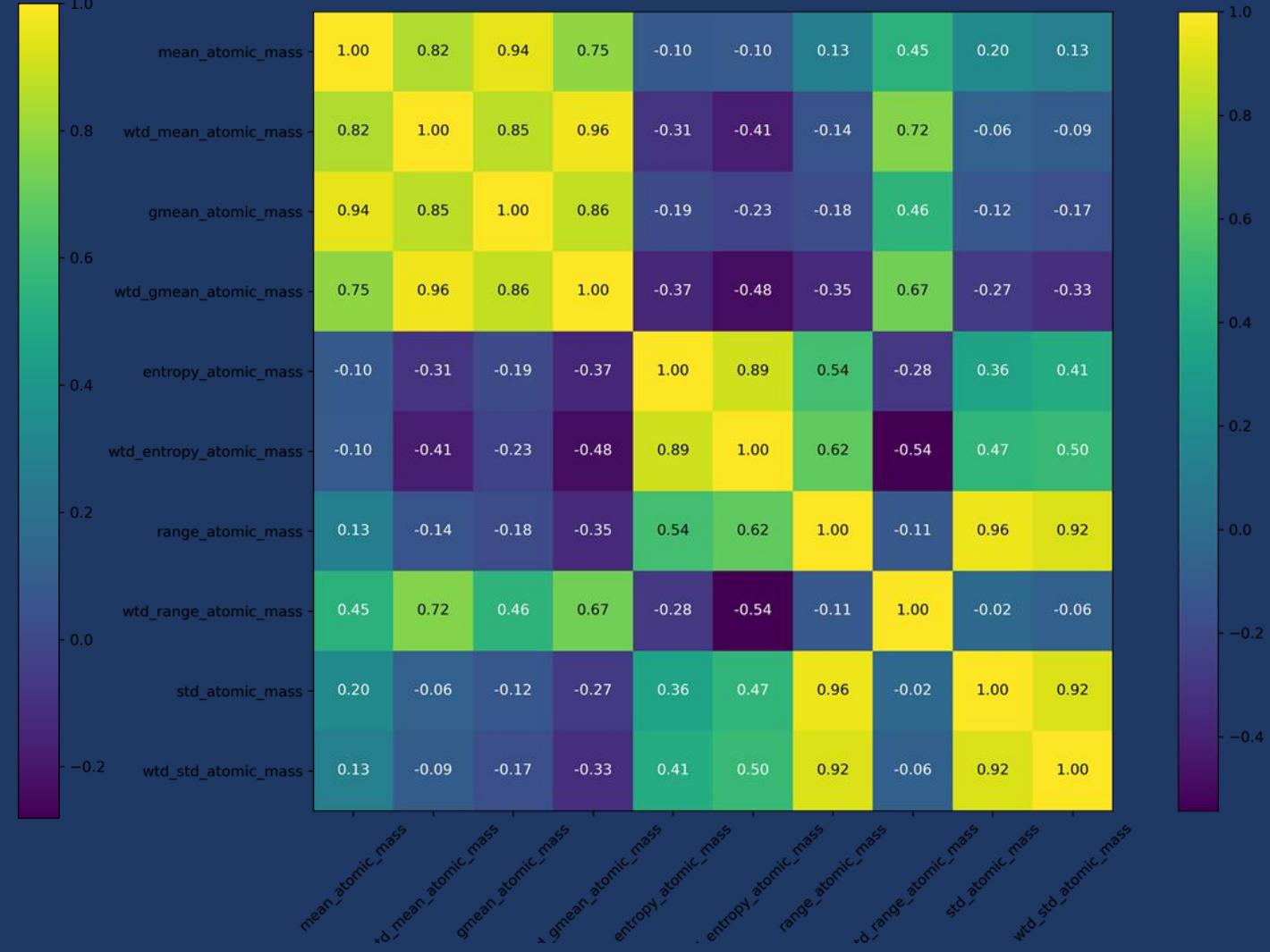
- It measures the linear dependence between pairs of features
- $-1 \leq r \leq 1$
- $r = 0$  : no correlation
- $r = 1$  : positive correlation
- $r = -1$  : negative correlation

# Relationship between the statistics of the same physical property

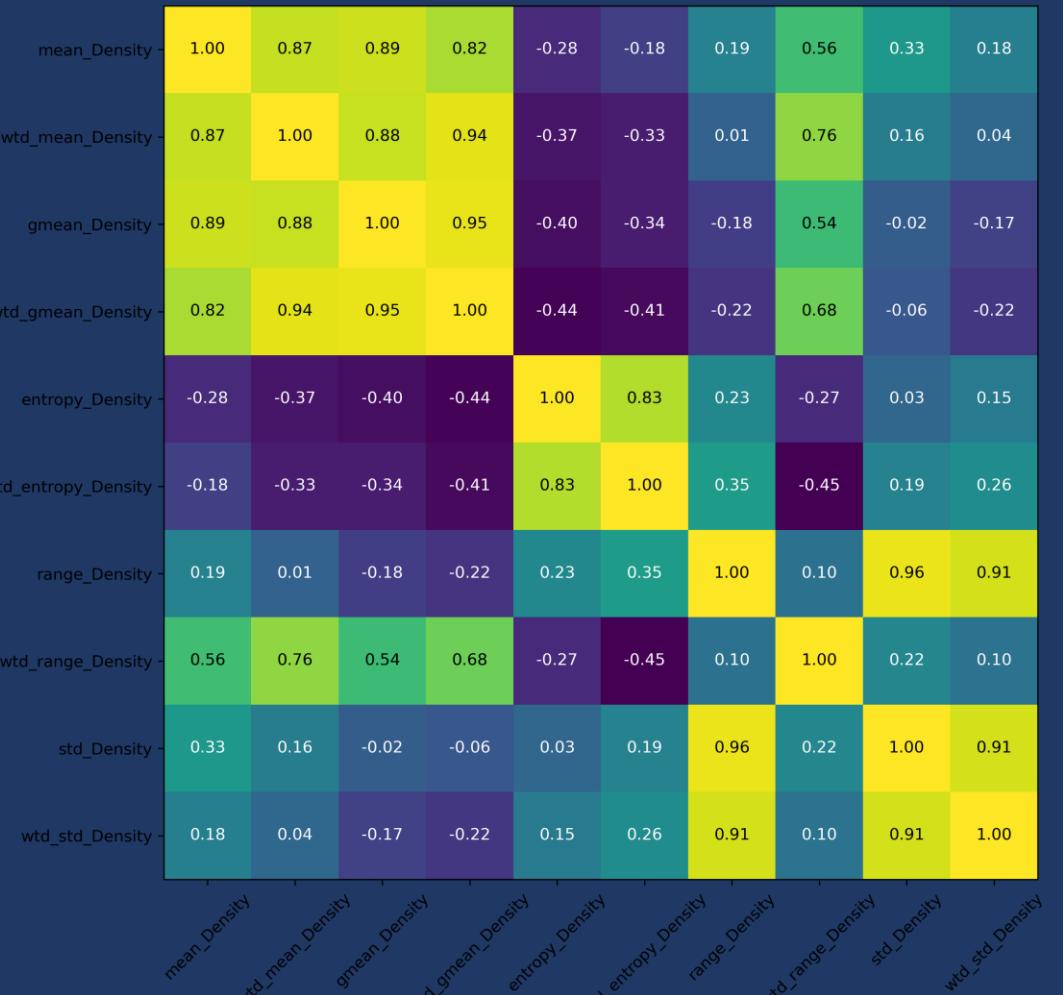
First Ionization Energy



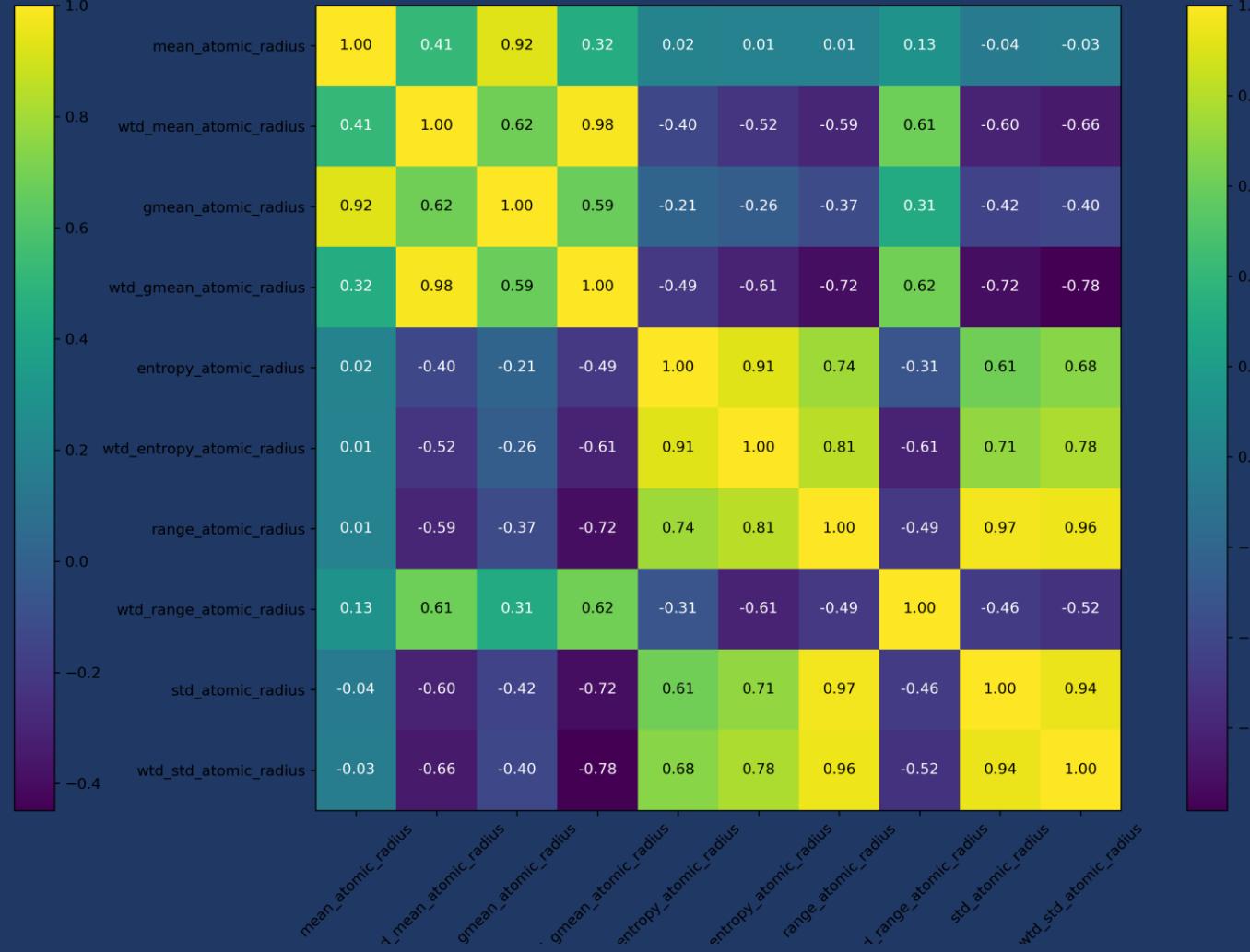
Atomic Mass



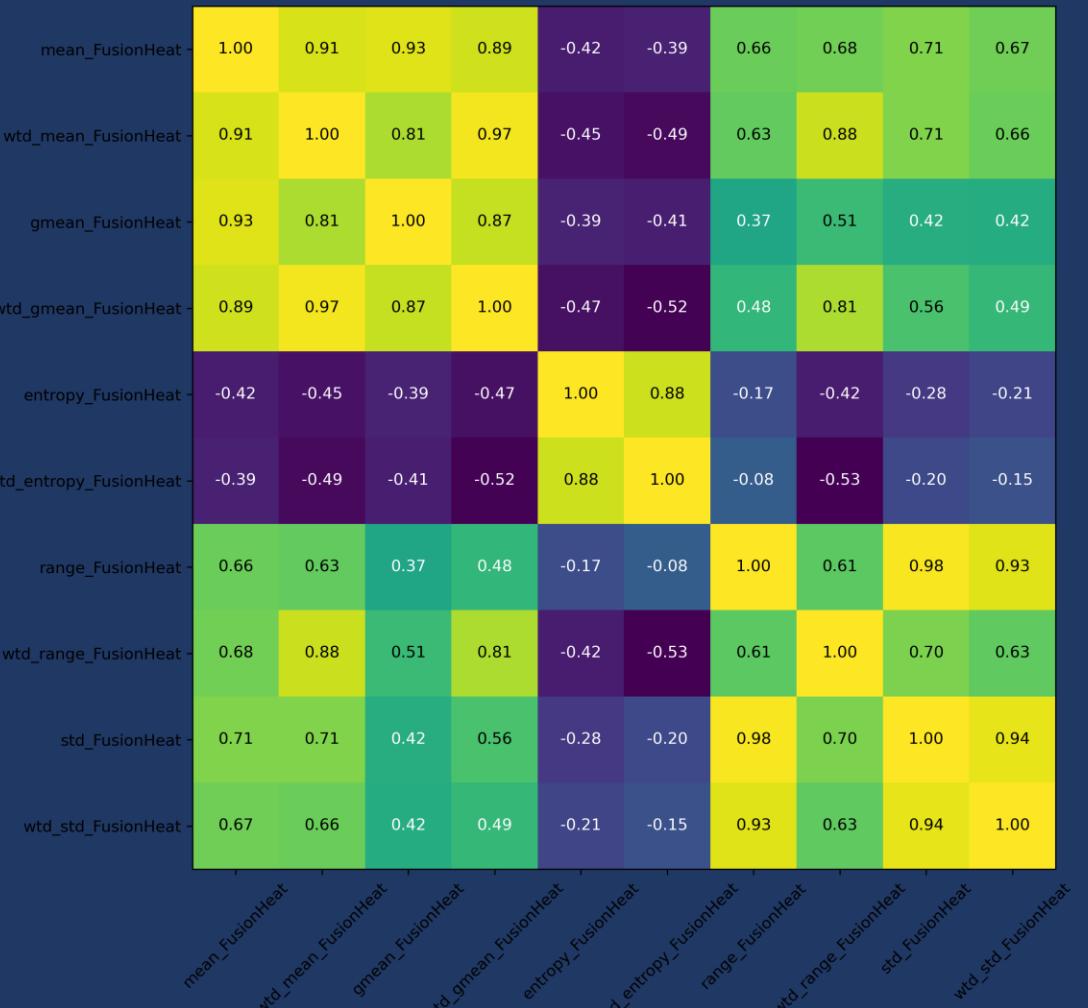
Density



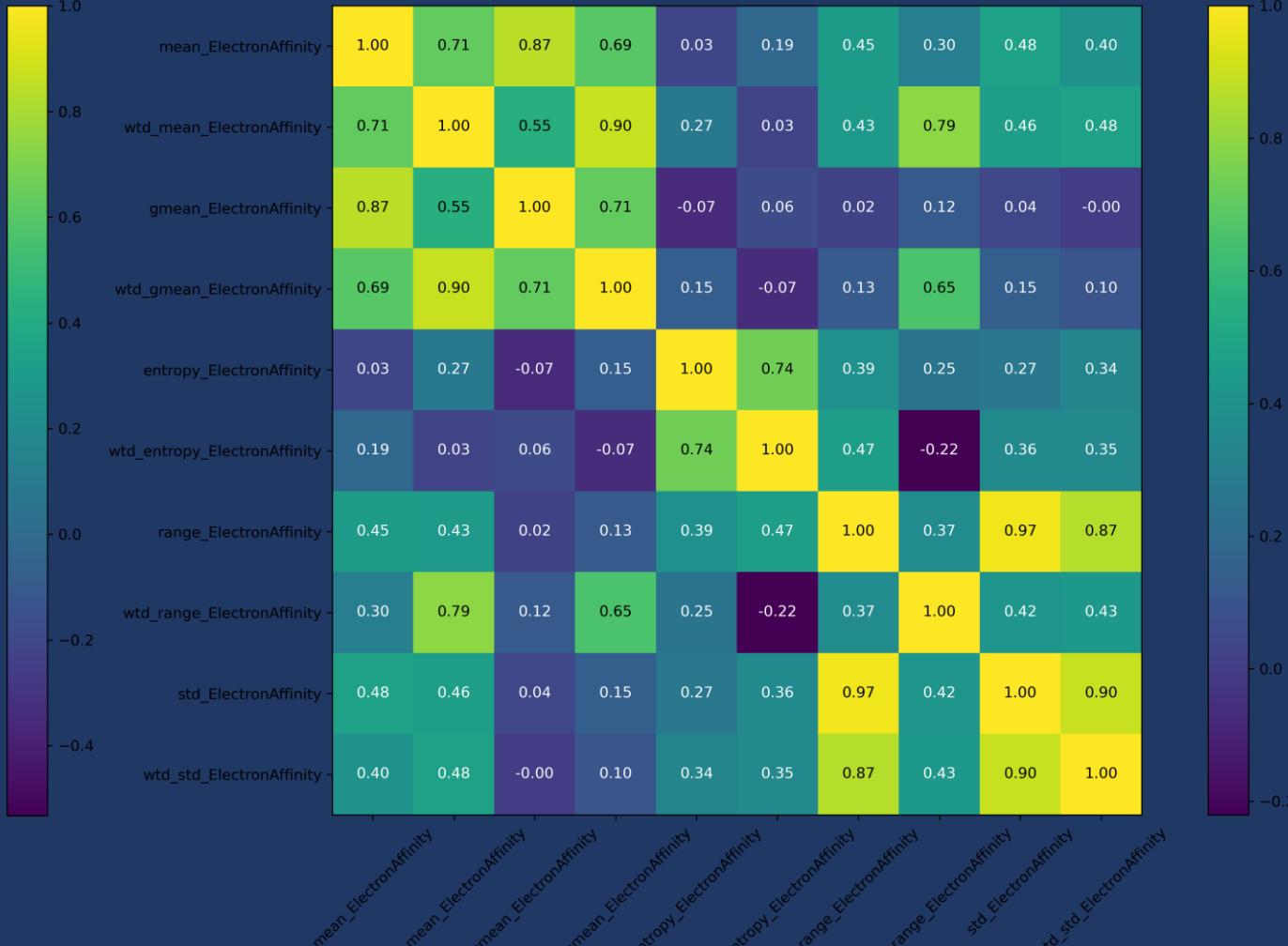
Atomic Radius



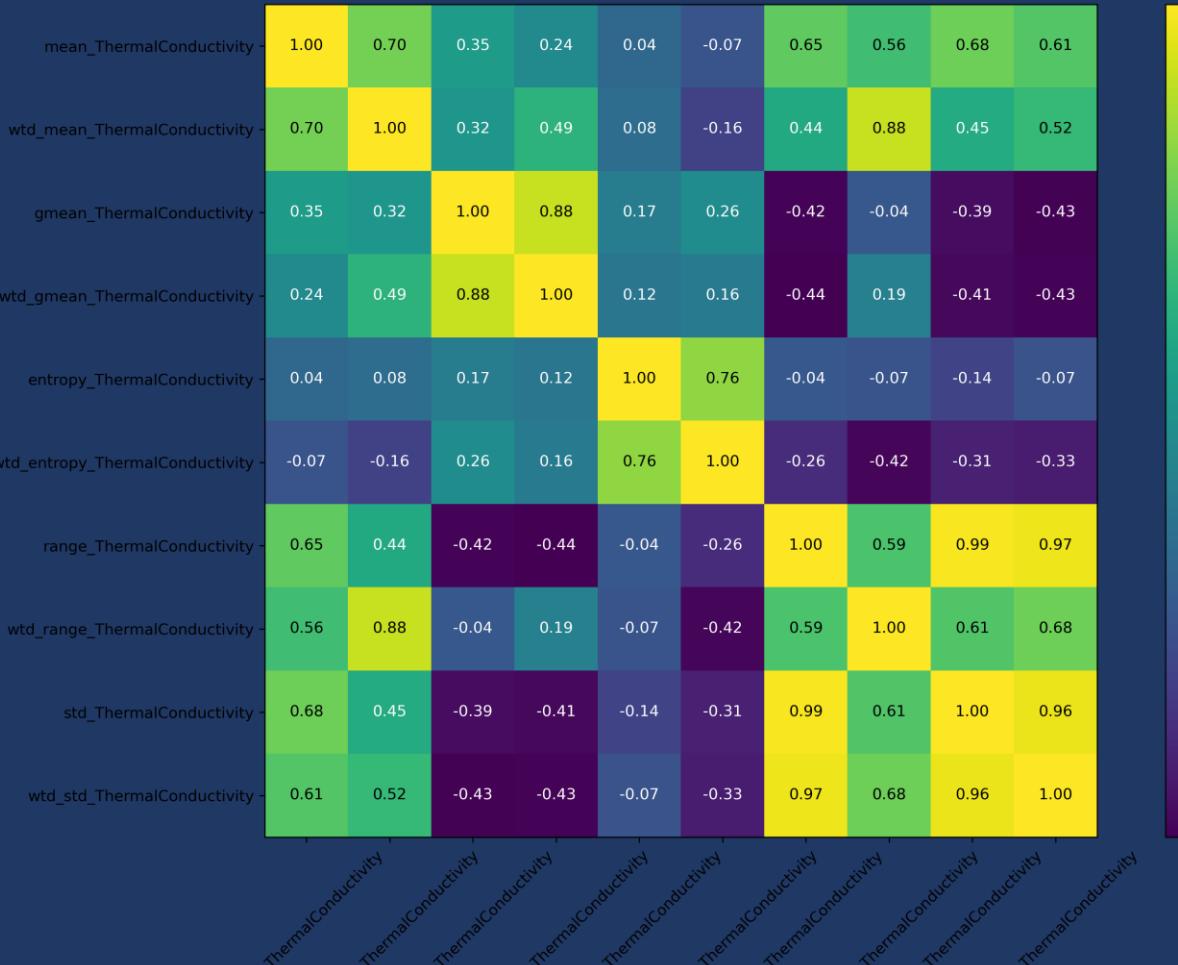
Fusion Heat



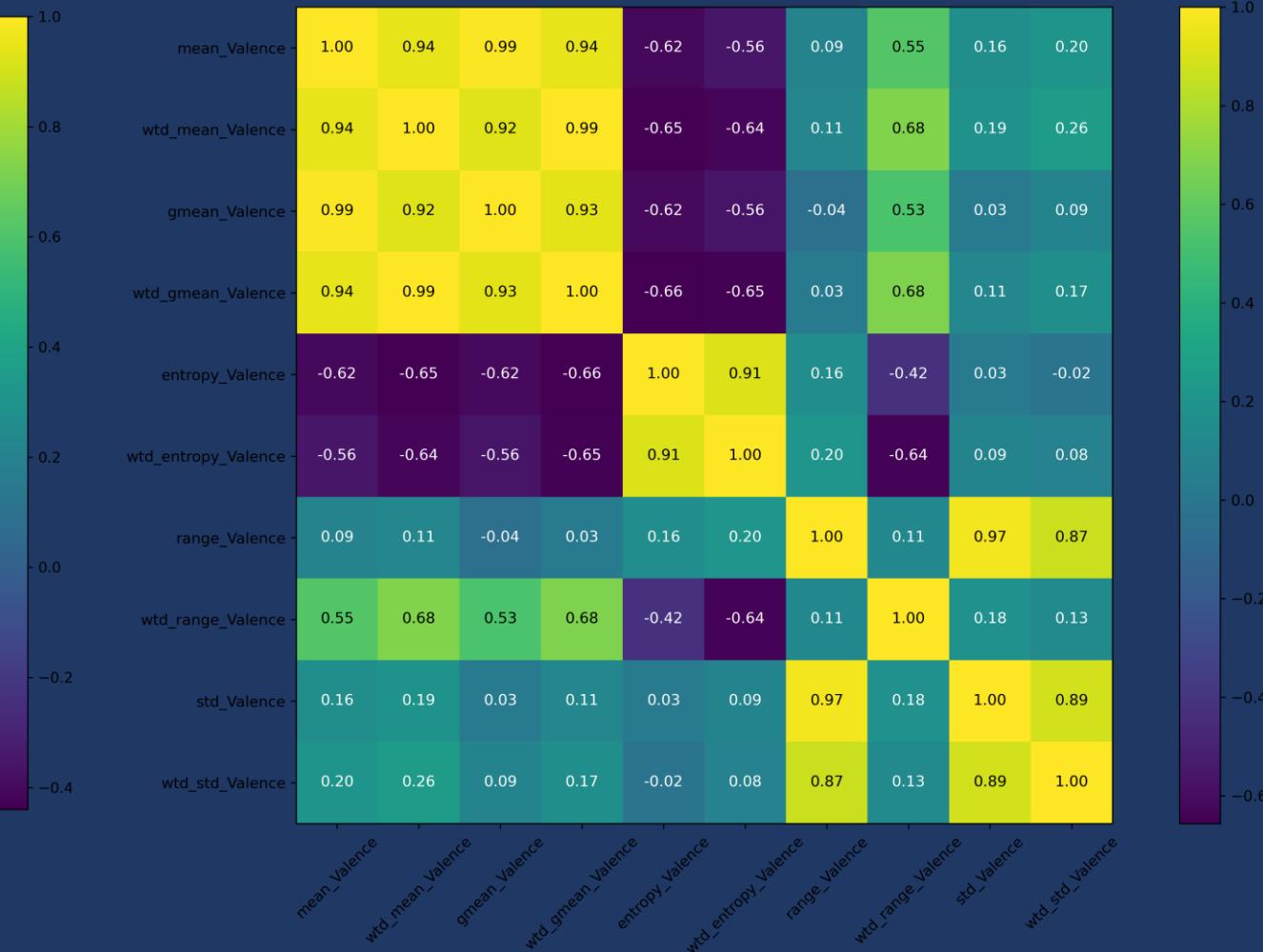
Electron Affinity



# Thermal Conductivity



# Valence

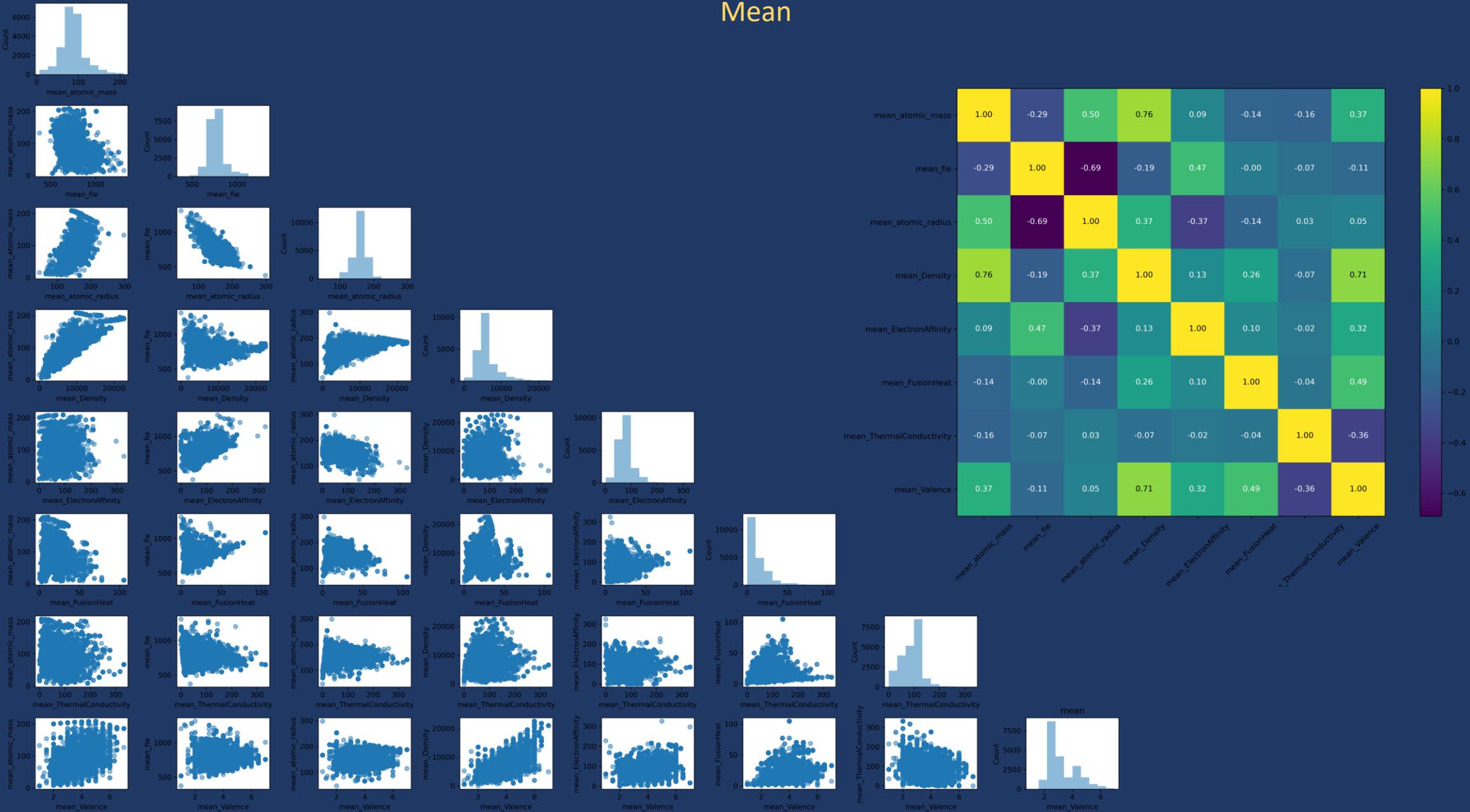


# Scatter plot Matrix

- Given a set of variables  $x_1, x_2, \dots, x_n$ , it contains all the pairwise scatter plots of the variables on a single page in a matrix format. The i-th row and j-th column of this matrix is a plot of  $x_i$  vs  $x_j$ .
- Since  $x_i$  vs  $x_j$  is equivalent to  $x_j$  vs  $x_i$  with the axes reversed, we prefer to omit the plots above the diagonal.
- On the diagonal, instead of graphing a trivial relationship, we plot the univariate histogram.

Now we will explore the relationship between feature that represents different physical properties, fixing a single statistic.

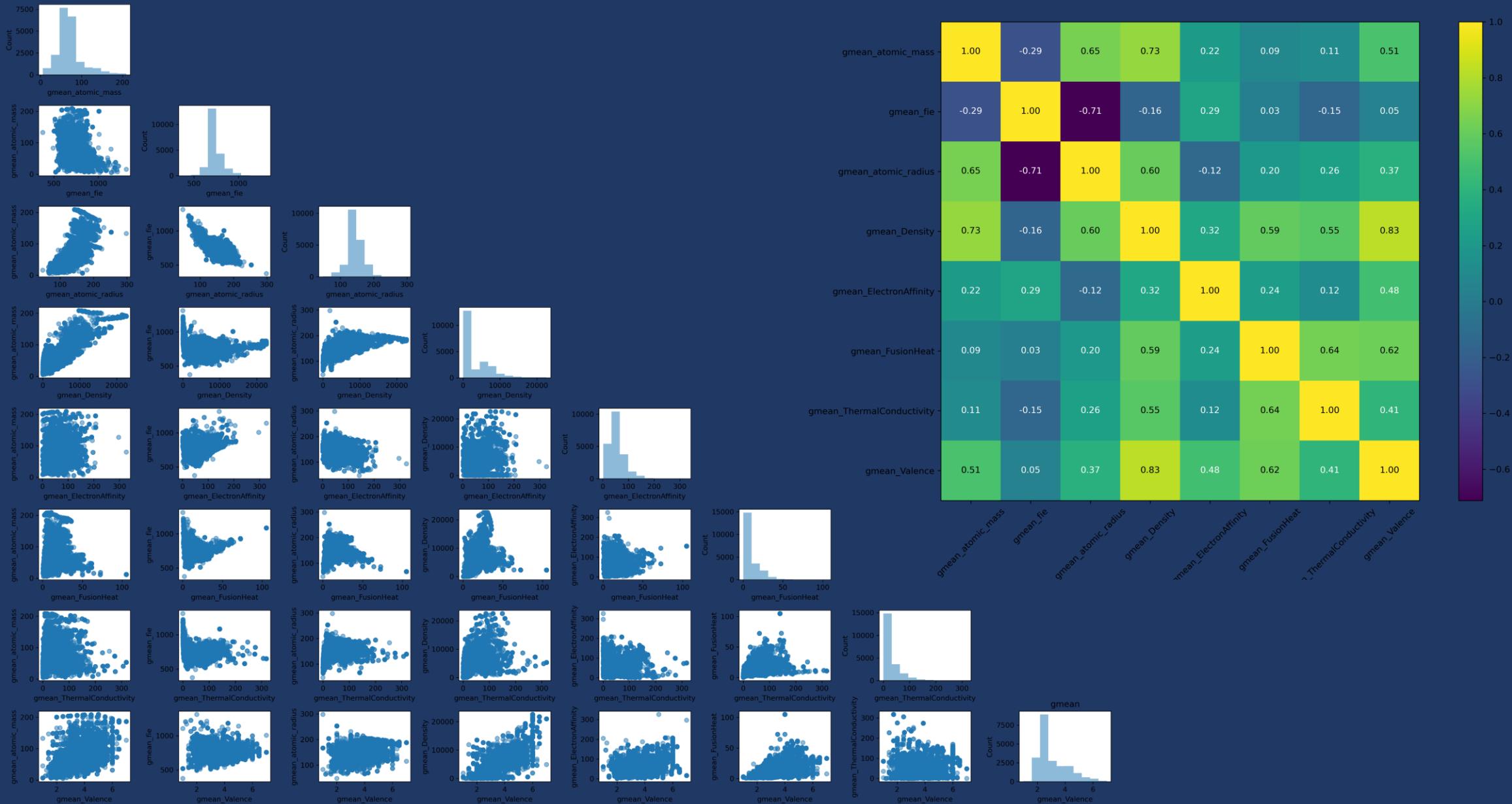
## Mean



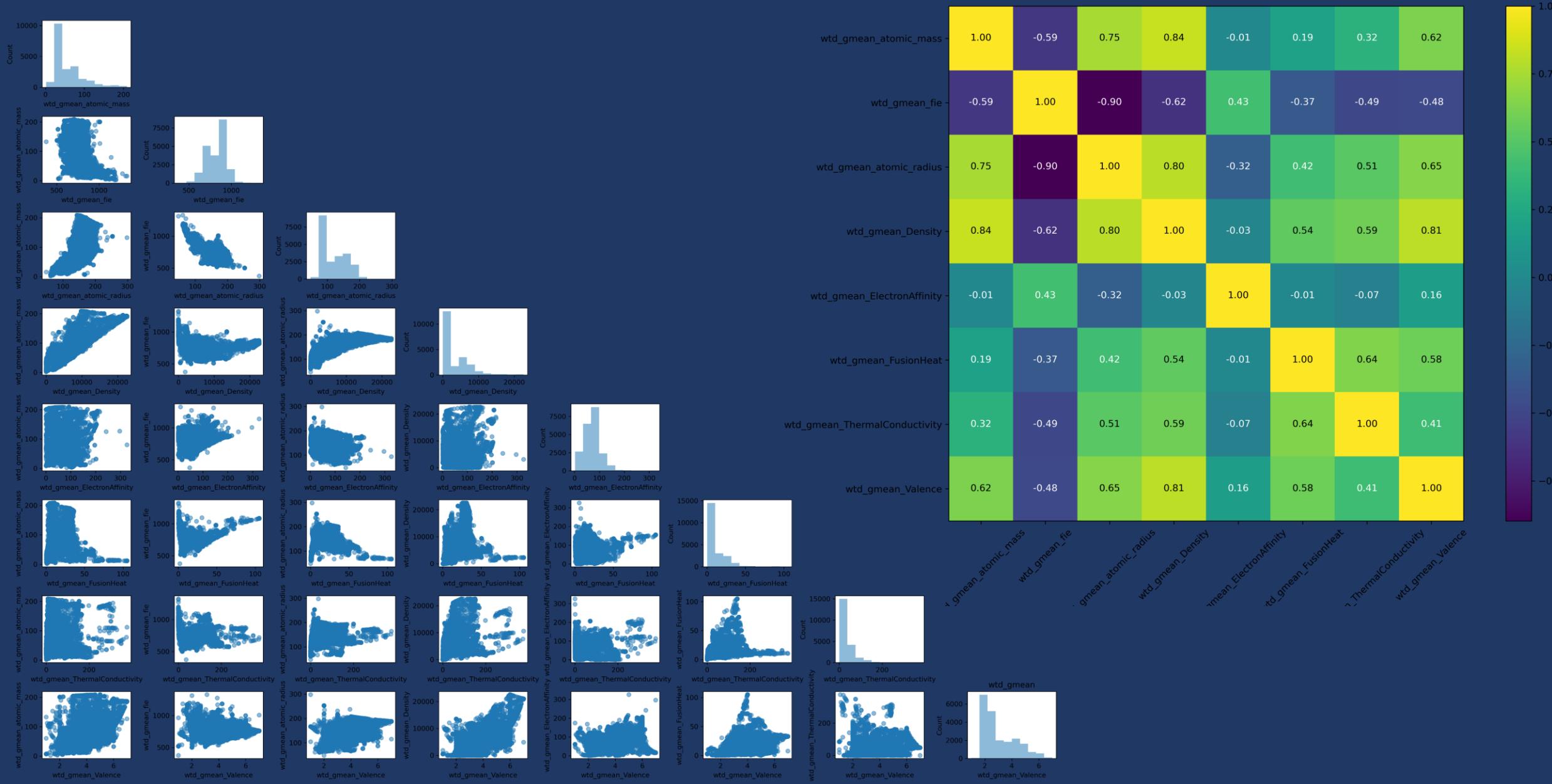
## Weighted Mean



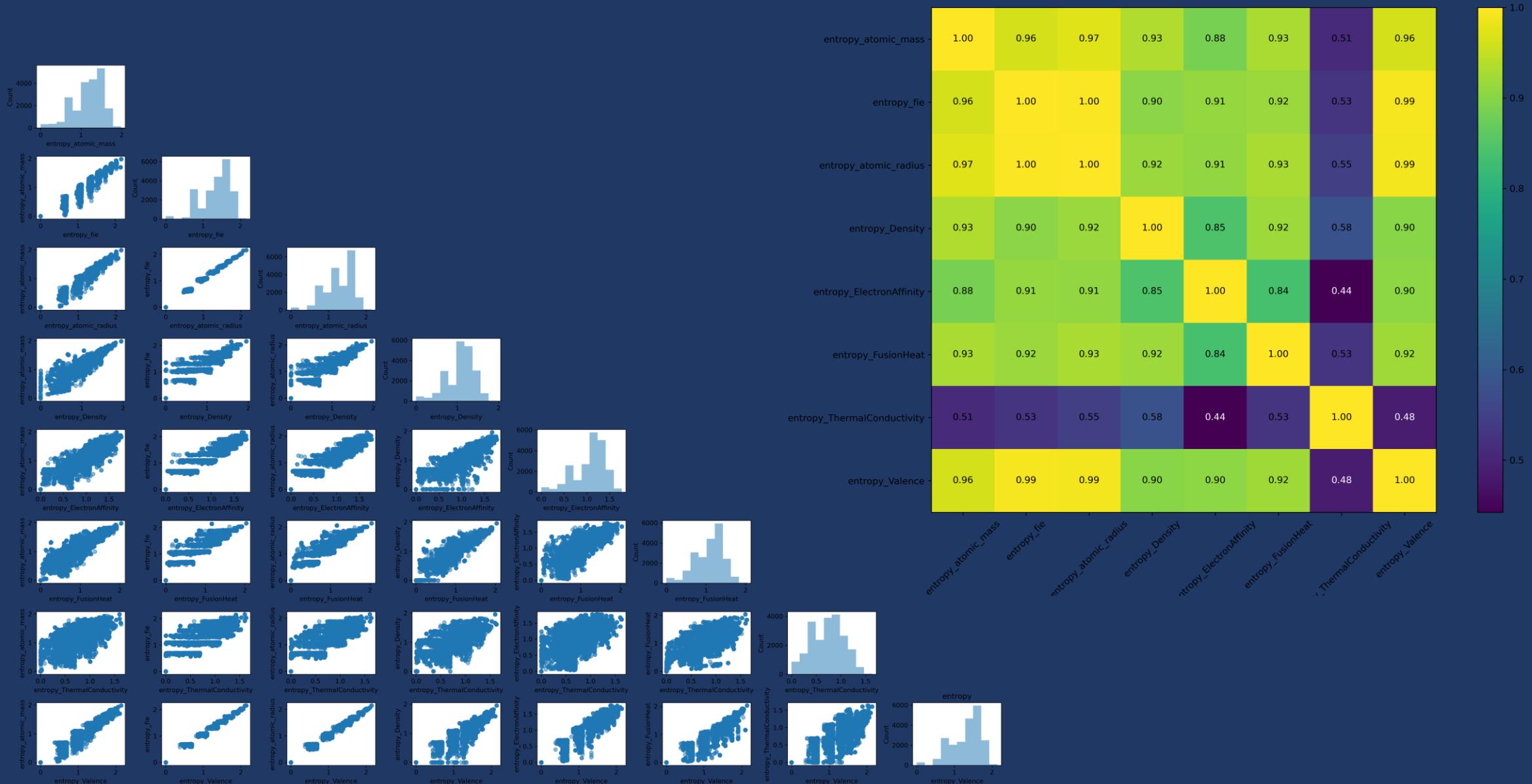
## Geometric Mean



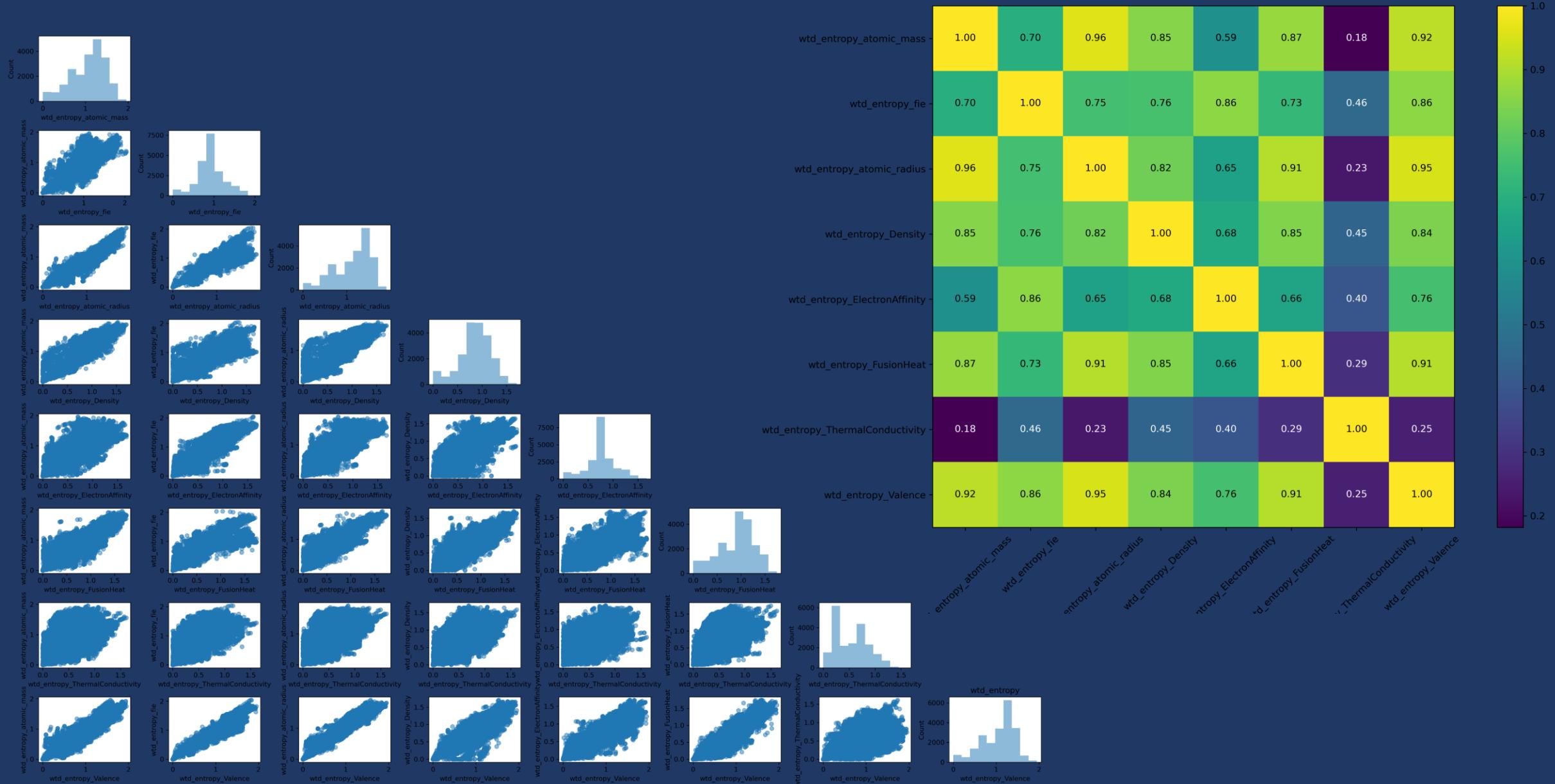
# Weighted Geometric Mean



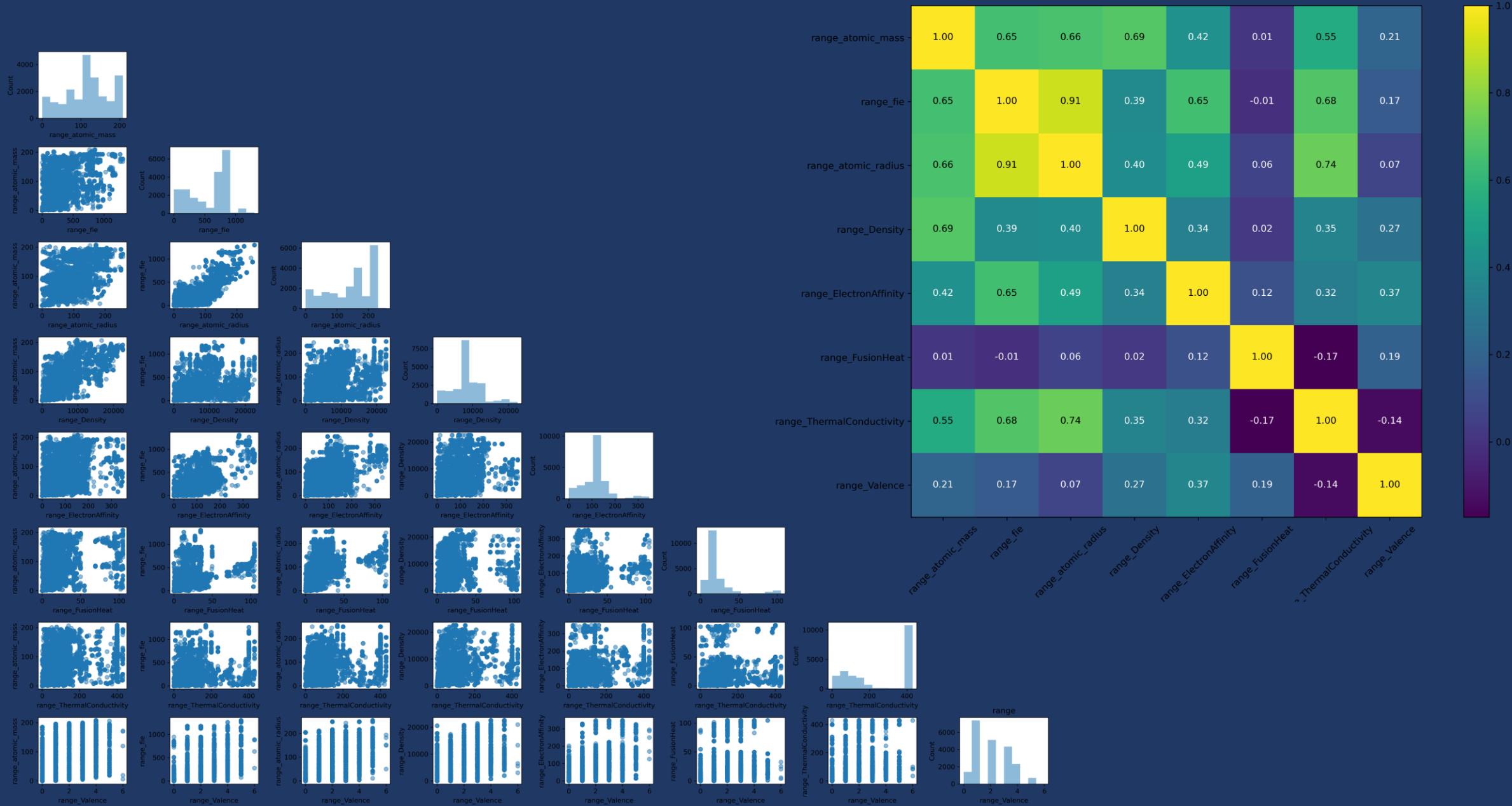
# Entropy



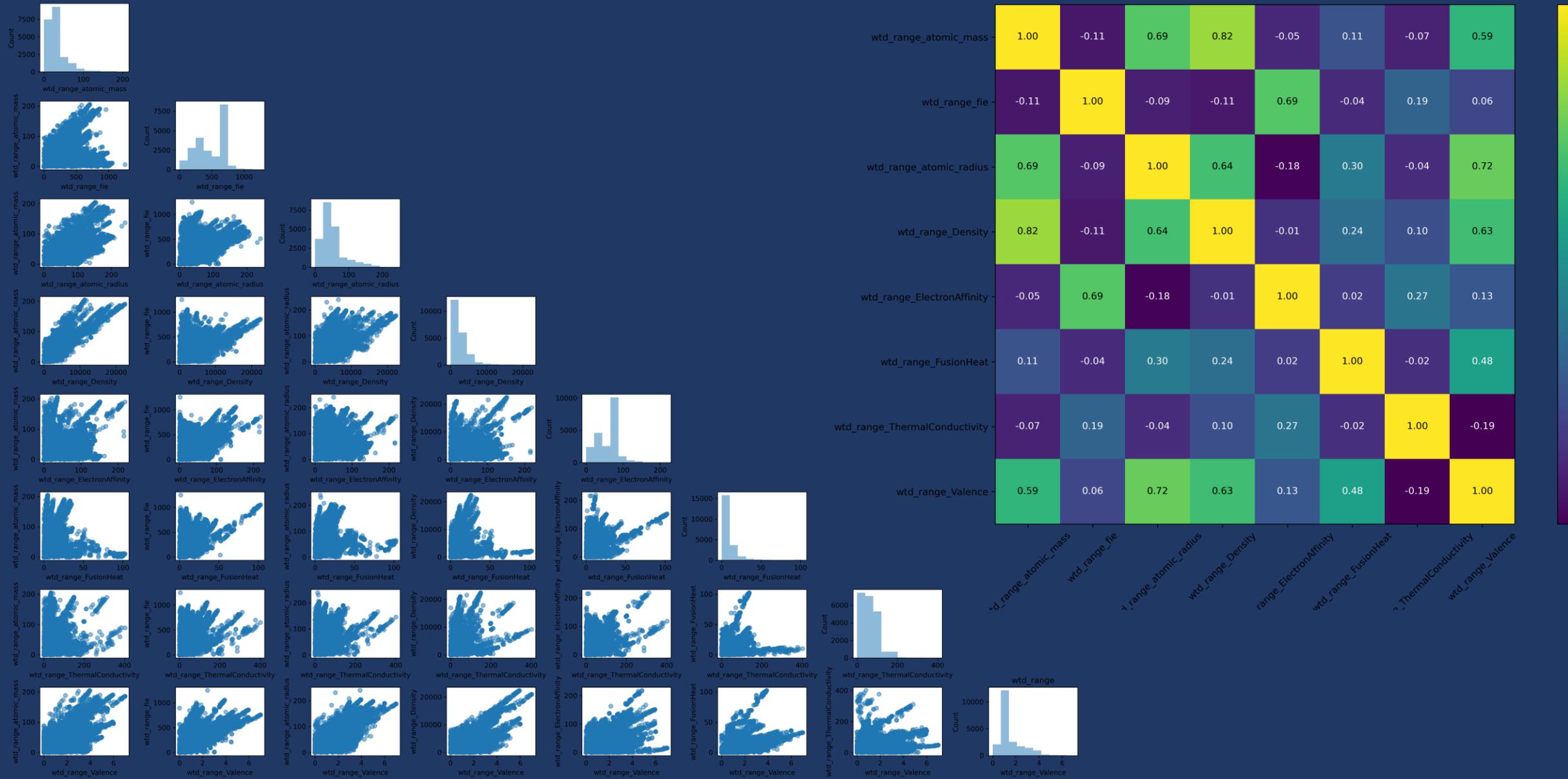
# Weighted Entropy



# Range



# Weighted Range



# Standard Deviation



# Weighted Standard Deviation



# Feature Scaling

## PROBLEMS

- Distorted feature importances on some models
- Overfitting
- Long Running time

## SOLUTIONS

- Feature Standardization:

$$x' = \frac{x - \mu}{\sigma}$$

- Feature Normalization:

$$x' = a + \frac{(x - \min(x))(b - a)}{\max(x) - \min(x)}$$

# Dimensionality Reduction

## PROBLEMS

- Multicolliniarity
- Overfitting
- High Computational Cost

## SOLUTIONS

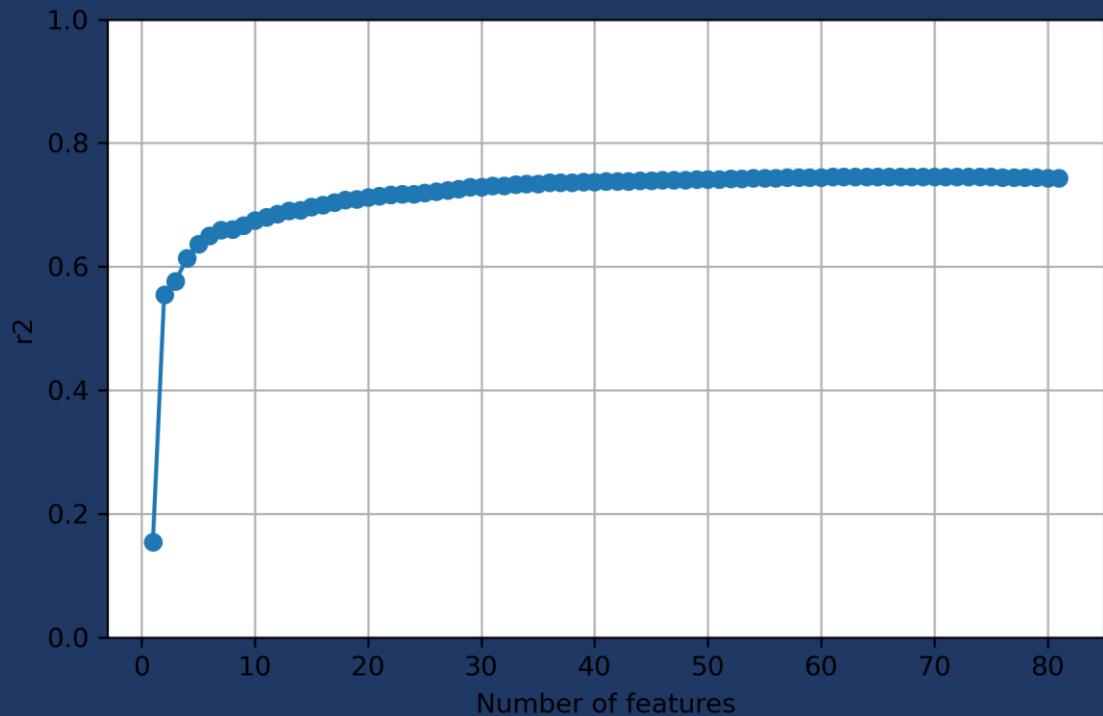
- Feature selection
- Feature extraction

# Feature Selection

- Variance Threshold Selection
  - It removes all features whose variance is lower than a certain threshold.
  
- Sequential Backward Selection
  - Gridy Search algorithm
  - Useful for algorithms that don't support regularization
  
- Random Forest
  - Feature importances

# Feature Selection on Superconductors Dataset

SBS on X\_train\_std



30 features selected

Dropping low variance features

Variance Threshold	Feature selected
10	55

# Feature Extraction

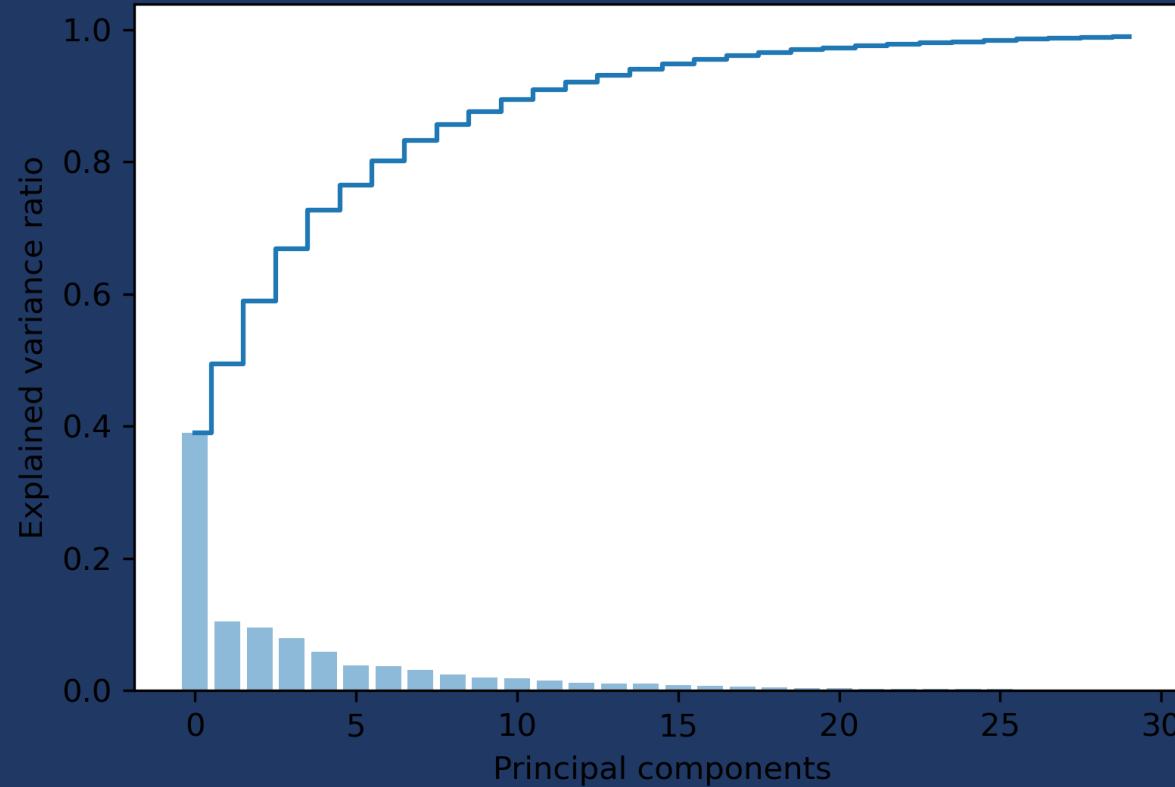
- Feature extraction transform or project the data onto a new feature space of lower dimensionality than the original one.
- Using feature selection algorithms we maintain the original features, but on the contrary using feature selection we create new features.
- The best-known example is the Principal Component Analysis algorithm.

# Principal Component Analysis

PCA can be summarised in the following steps:

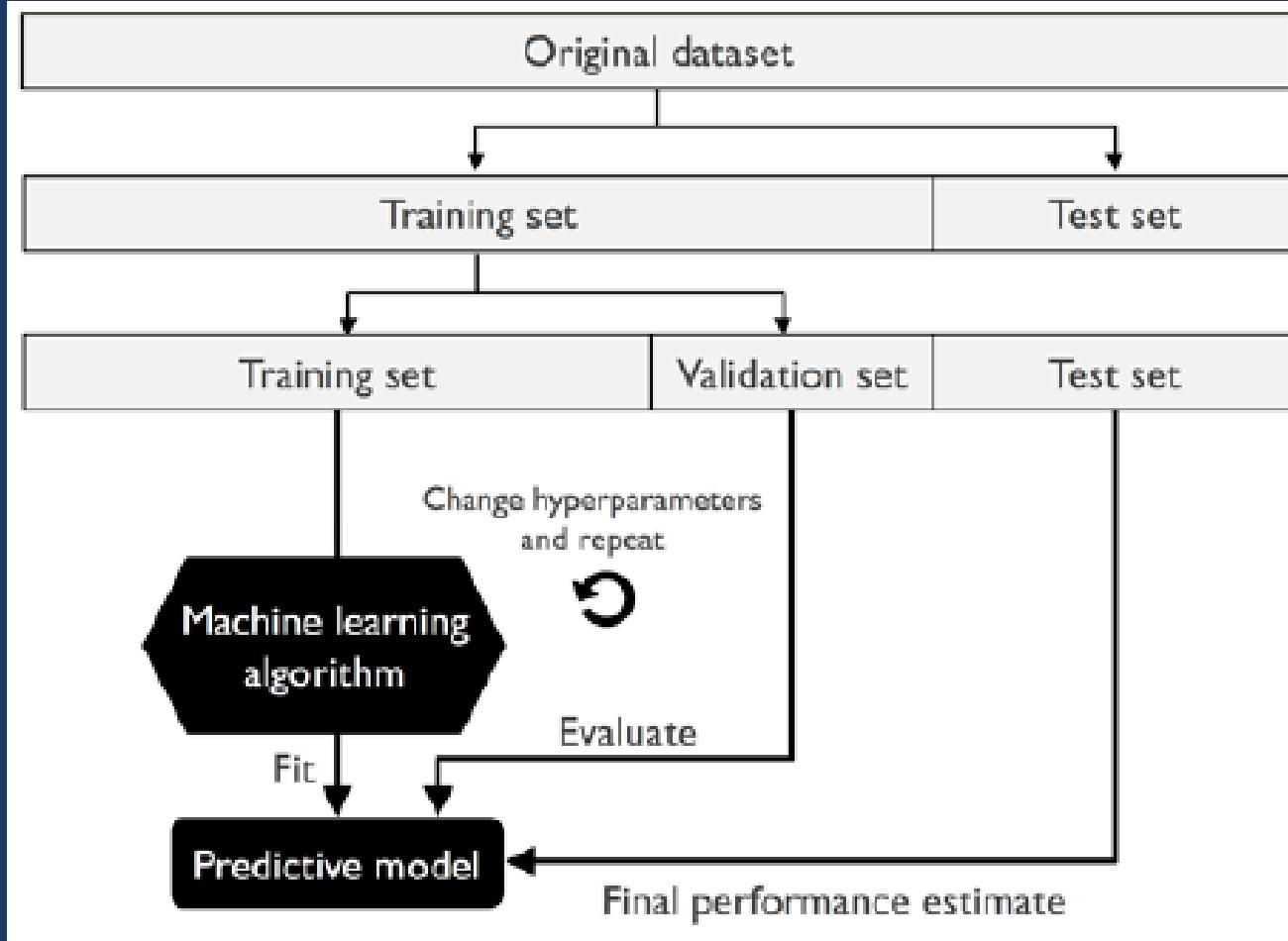
1. Standardize the d-dimensional dataset.
2. Construct the covariance matrix.
3. Decompose the covariance matrix into its eigenvectors and eigenvalues.
4. Sort the eigenvalues by decreasing order to rank the corresponding eigenvectors.
5. Select k eigenvectors which correspond to the k largest eigenvalues, where k is the dimensionality of the new feature subspace( $k < d$ ).
6. Construct a projection matrix W from the "top" k eigenvectors.
7. Transform the d-dimensional input dataset X using the projection matrix W to obtain the new k-dimensional feature subspace

# PCA on Superconductors Dataset



- 30 features selected.
- Explained variance ratio with 30 components: 0,99.

# Model Performance



- PARAMETER VS HYPERPARAMETER

# Cross Validation

It estimates the generalization performance of machine learning models

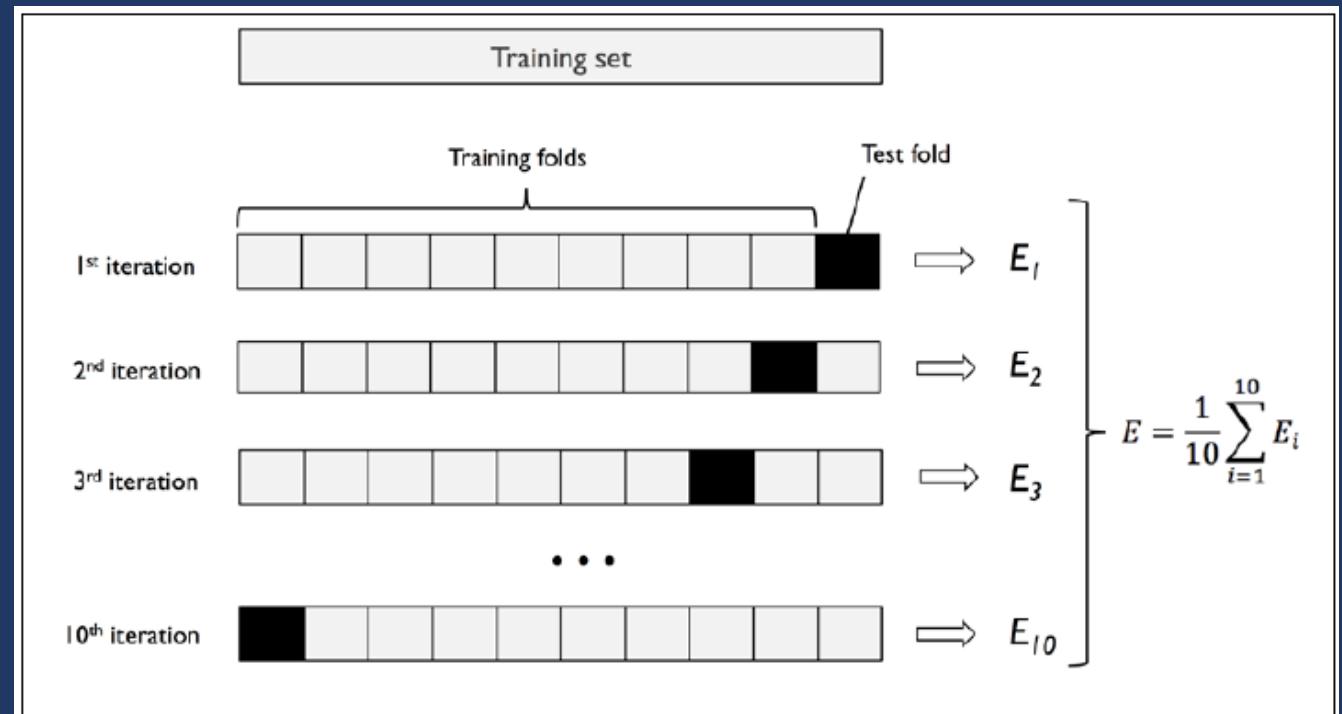
CROSS VALIDATION ALGORITHM	DESCRIPTION
HOLDOUT	<ul style="list-style-type: none"><li>• It split the initial dataset into separate training and test datasets or better into 3 parts: training, validation and test datasets.</li></ul>
K - FOLD	<ul style="list-style-type: none"><li>• The performance estimate is less sensitive to the sub-partitioning of the training data.</li><li>• It randomly split the training dataset into <math>k</math> folds without replacement, where <math>k - 1</math> folds are used for the model training, and one fold is used for performance evaluation.</li></ul>
LEAVE - p - OUT	<ul style="list-style-type: none"><li>• It's a k-fold CV generalization, where you choose the number <math>p</math> of folds used for the model training, Then, <math>n-p</math> folds are used for the evaluation.</li></ul>

# Cross Validation on Superconductors dataset

- HOLDOUT

- 20% Test Dataset
- 80% Training Dataset (of which 20% is the Validation dataset)

- 5-FOLD CV



# HYPERPARAETER TUNING

## Grid Search CV

- Brute-force exhaustive search paradigm
- Sklearn.model\_selection.GridSearchCV

### PARAMETERS:

- estimator: algorithm of which you want to optimize the hyperparameters
- param\_grid : collection of hyperparameters to be tested in the algorithm
- scoring: Strategy to evaluate the performance of the cross-validated model on the test set
- n\_jobs: Number of jobs to run in parallel
- refit: Refit an estimator using the best found parameters on the whole dataset
- n\_iter :Number of parameter settings that are sampled (only for Random Search CV)

# Validation Curves

- Validation curves are a useful tool for improving the performance of a model by addressing issues such as overfitting or underfitting.
- Validation curves are plots of the training and test scores as functions of the model parameters, such as the `max_depth` of decision tree.

# Evaluating the Performance of Regression Models

$$MSE = \frac{1}{n} \sum_{\{i=1\}}^n (y^{(i)} - \hat{y}^{(i)})^2$$

**Mean Squared Error**

$$SSE = \sum_{\{i=1\}}^n (y^{(i)} - \hat{y}^{(i)})^2$$

**Sum of Squares  
Error**

$$SST = \sum_{\{i=1\}}^n (y^{(i)} - \mu_y)^2$$

**Sum of Squares  
Total**

$$R^2 = 1 - \frac{SSE}{SST}$$

**Coefficient of Determination**

For training data:  $0 \leq R^2 \leq 1$ , but for test data it can become negative

# Supervised Learning Techniques for Regression Problem

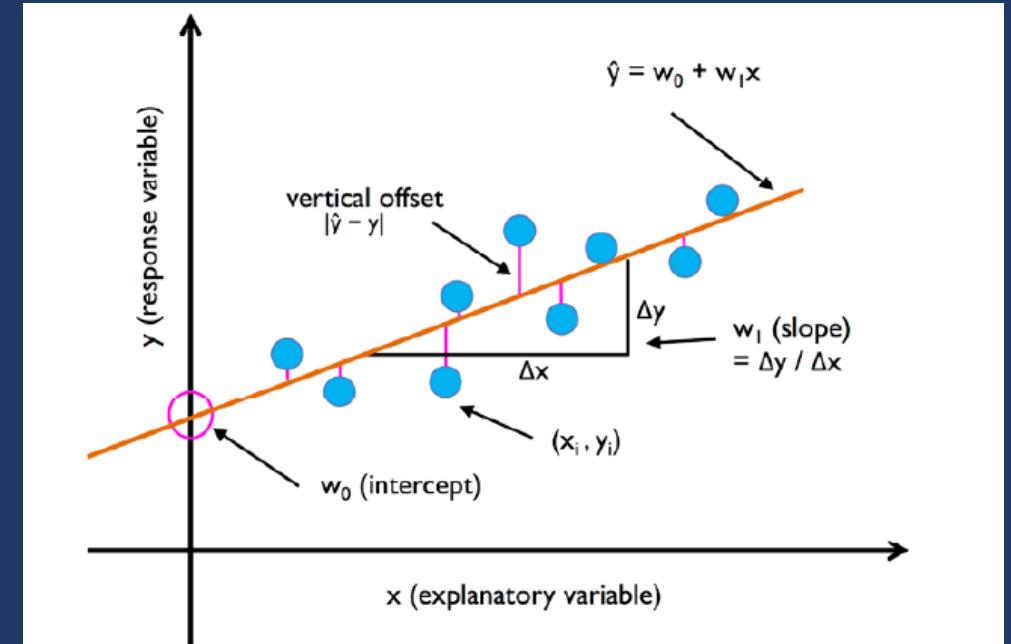
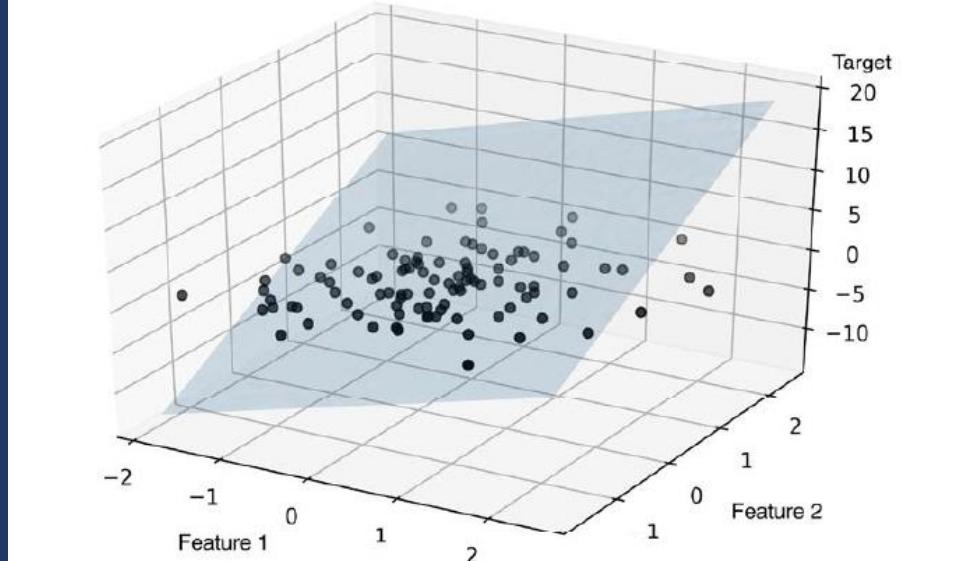


# Multiple Linear Regression

The goal of linear regression is to model relationship between *one or multiple* features and a **continuous** target variable.

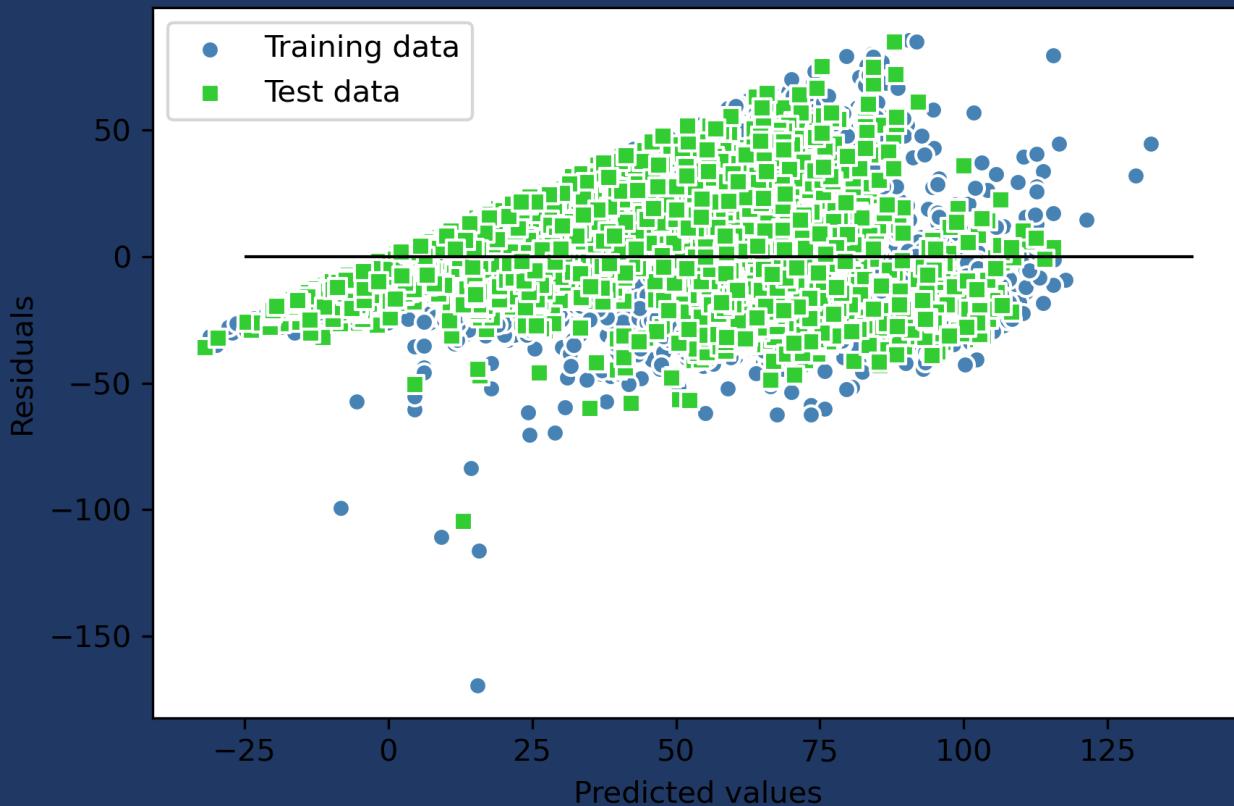
$$y = w_0 + w_1 x$$

$x$  = explanatory variable(s)  
 $y$  = response variable  
 $w_0$  = intercept  
 $w_1$  = slope(s)



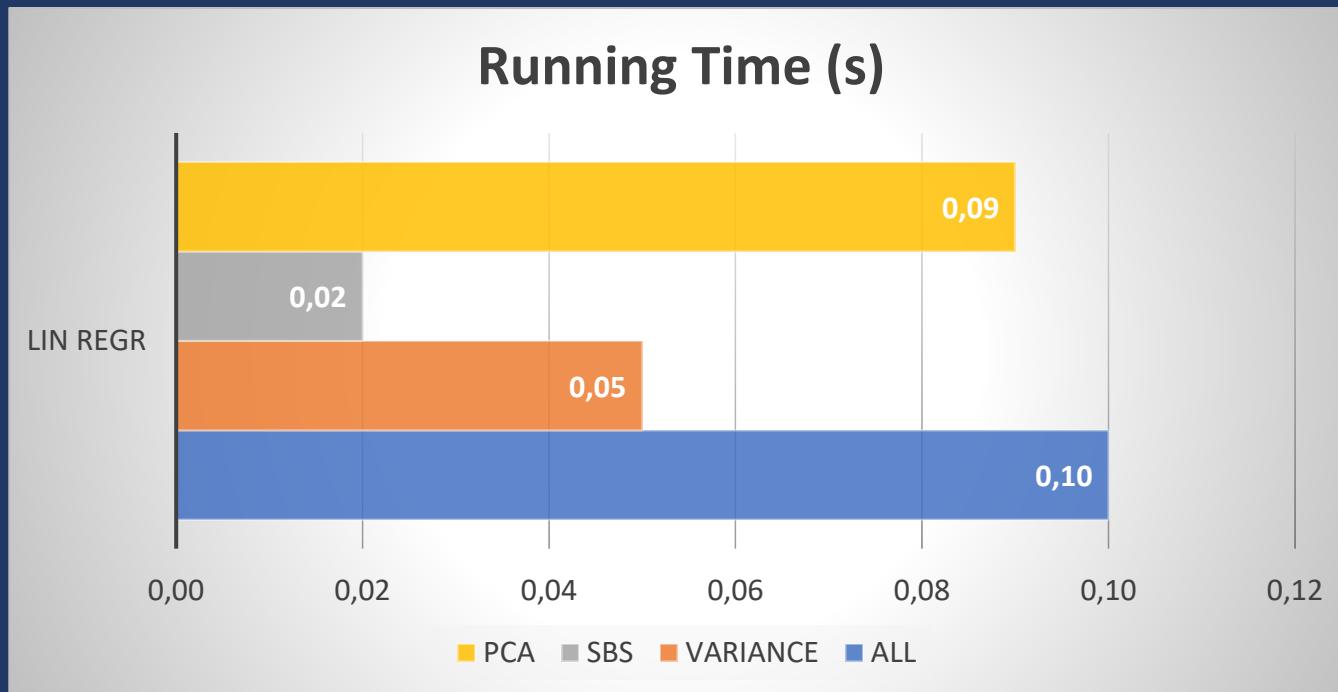
PROS	CONS
Simple and Rapid algorithm	Heavy impact by the presence of outliers

## RESIDUAL PLOT



- Residues follow a fairly regular growing pattern, indicating that the data is not "very linear"

METRICS	ALL	VARIANCE	SBS	PCA
R2 TRAIN	0,74	0,72	0,720	0,74
R2 TEST	0,73	0,710	0,710	0,73



# Regularized Methods for Regression

## Lasso

- L1 penalized model
- $J(w)_{LASSO} = \sum_{\{i=1\}}^n (y^{(i)} - \hat{y}^{(i)})^2 + \lambda |w|_1$
- Parameter: 'alpha'

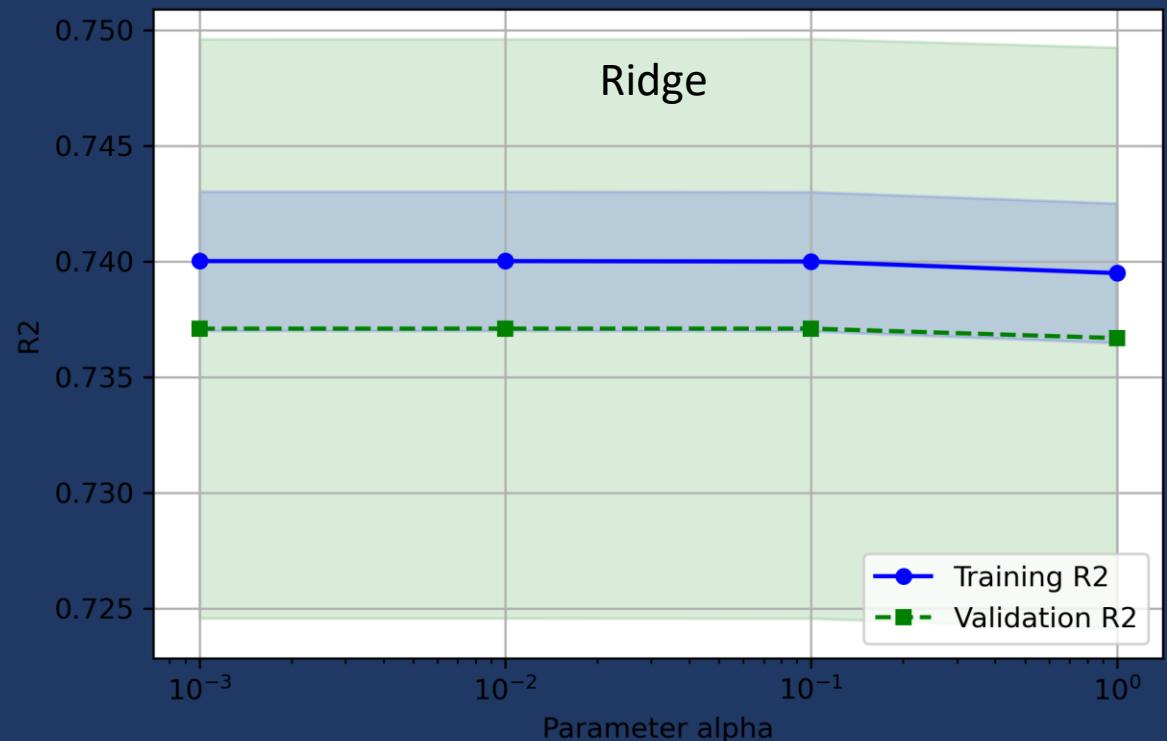
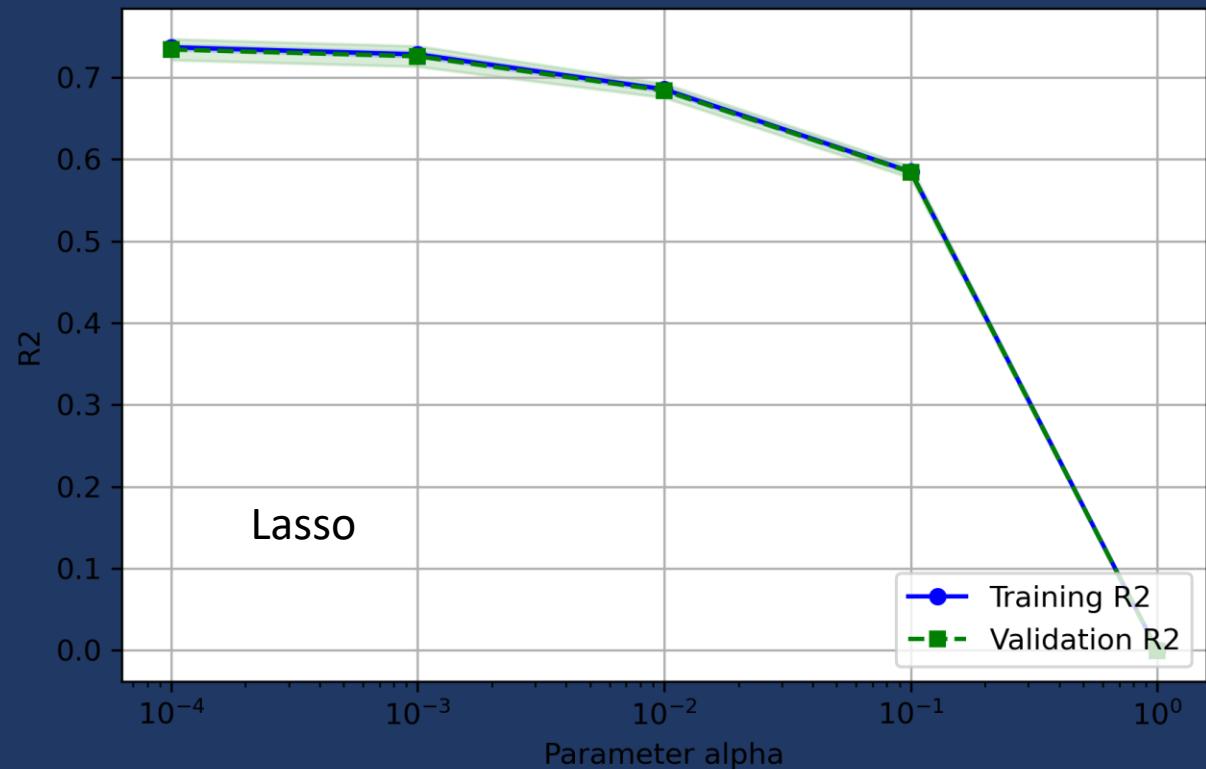
## Ridge

- L2 penalized model
- $J(w)_{RIDGE} = \sum_{\{i=1\}}^n (y^{(i)} - \hat{y}^{(i)})^2 + \lambda |w|_2^2$
- Parameter: 'alpha'

## ElasticNet

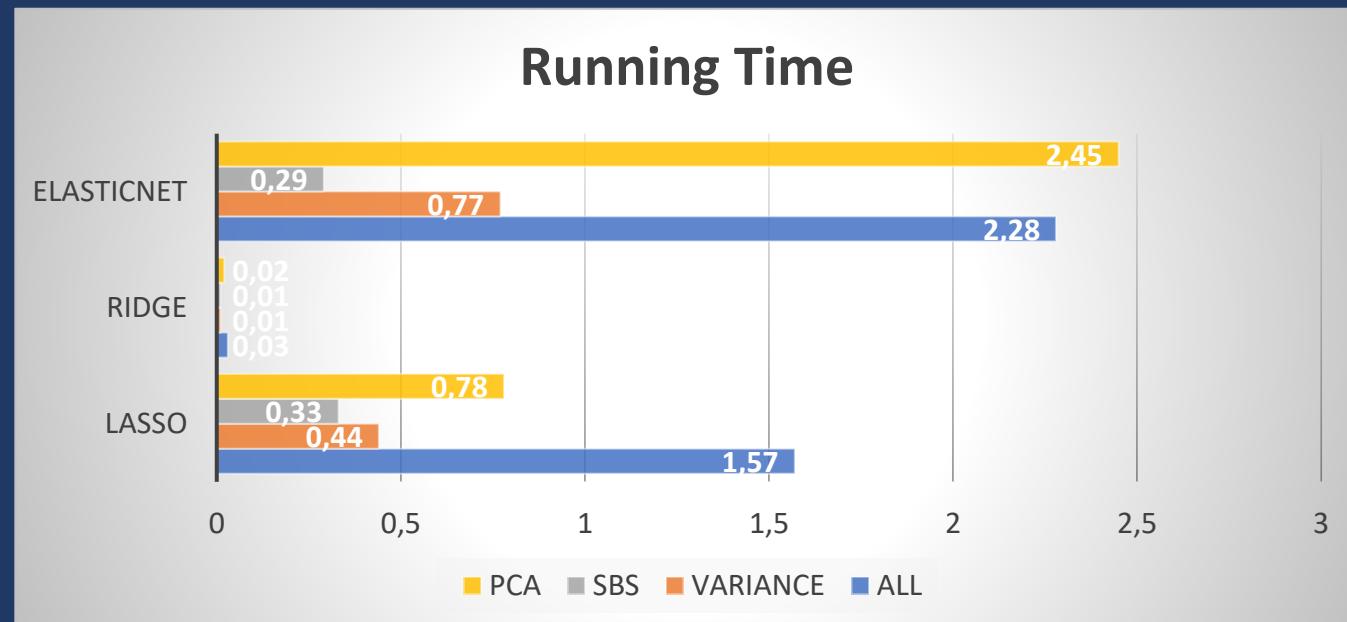
- Compromise between the two previous models
- $J(w)_{ELNET} = \sum_{\{i=1\}}^n (y^{(i)} - \hat{y}^{(i)})^2 + \lambda_1 \sum_{j=1}^m w_j^2 + \lambda_2 \sum_{j=1} |w_j|$
- Parameters: 'alpha', 'l1\_ratio'

# Tuning



ELASTICNET BEST PARAMETERS: 'l1'=0,001, 'alpha'=0,001

	METRICS	ALL	VARIANCE	SBS	PCA
LASSO	R2 TRAIN	0,62	0,63	0,62	0,62
	R2 TEST	0,62	0,62	0,62	0,62
RIDGE	R2 TRAIN	0,74	0,72	0,72	0,74
	R2 TEST	0,73	0,71	0,71	0,73
ELASTICNET	R2 TRAIN	0,61	0,58	0,59	0,61
	R2 TEST	0,61	0,58	0,59	0,61



# RANdom SAmple Consensus

## ALGORITHM

- 1) Select a random number of examples to be inliers and fit the model.
- 2) Test all other data points against the fitted model and add those points that fall within a user-given tolerance to the inliers.
- 3) Refit the model using all inliers.
- 4) Estimate the error of the fitted model versus the inliers.
- 5) Terminate the algorithm if the performance meets a certain user-defined threshold or if a fixed number of iterations were reached; go back to step 1 otherwise.

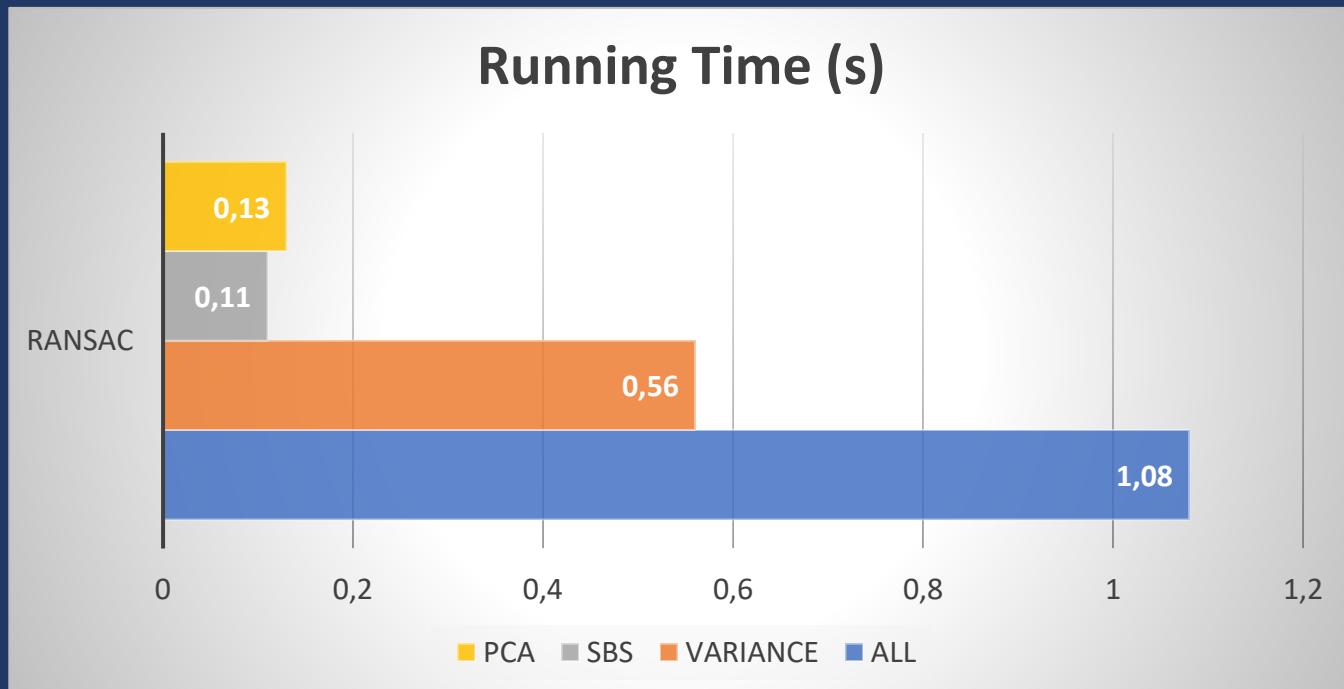
## PARAMETERS

- **min\_samples**: Minimum number of samples chosen randomly from original data
- **residual\_threshold**: Maximum residual for a data sample to be classified as an inlier
- **max\_trials**: Maximum number of iterations for random sample selection
- **loss**: ‘absolute\_error’ and ‘squared\_error’ are supported which find the absolute error and squared error per sample respectively

BEST PARAMETERS: `min_samples=1000, residual_threshold = 100, max_trials = 1000, loss= 'squared_error'`.

PROS	CONS
Reduce the potential effect of outliers	There are many hyperparameters to set

METRICS	ALL	VARIANCE	SBS	PCA
R2 TRAIN	0,72	0,71	0,71	0,74
R2 TEST	0,71	0,69	0,70	0,73



RANSAC

# Decision Tree Regressor

## INFORMATION GAIN

$$IG(D_p) = I(D_p) - \sum_{j=1}^m \frac{N_j}{N_p} I(D_j)$$

## IMPURITY METRIC

- Mean Squared Error

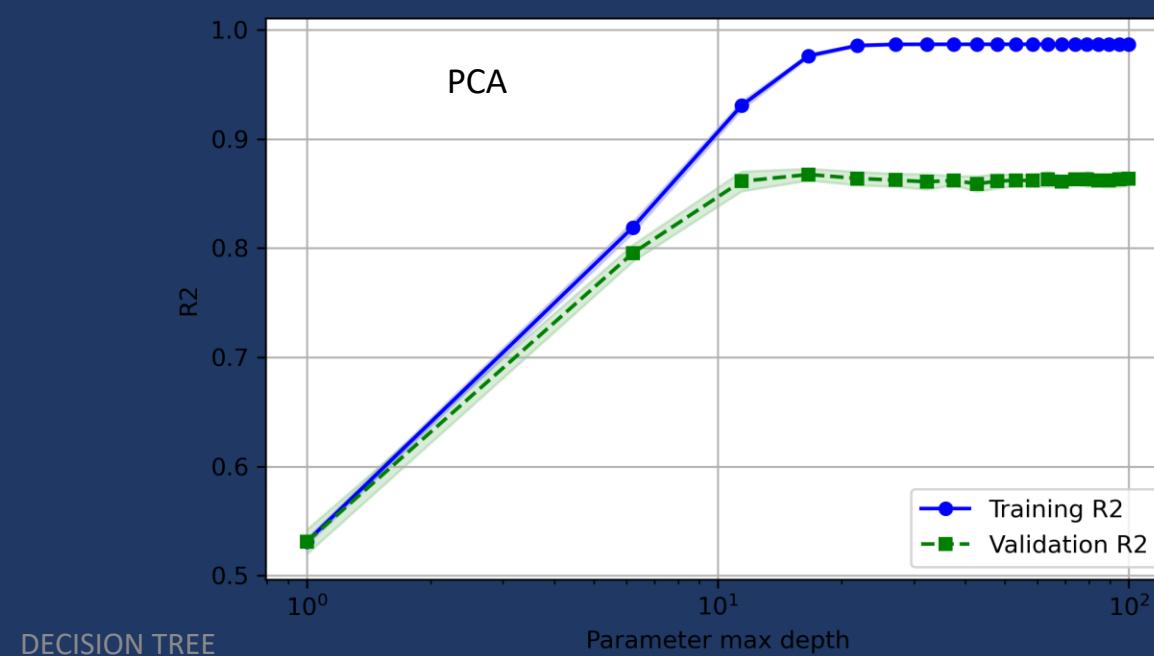
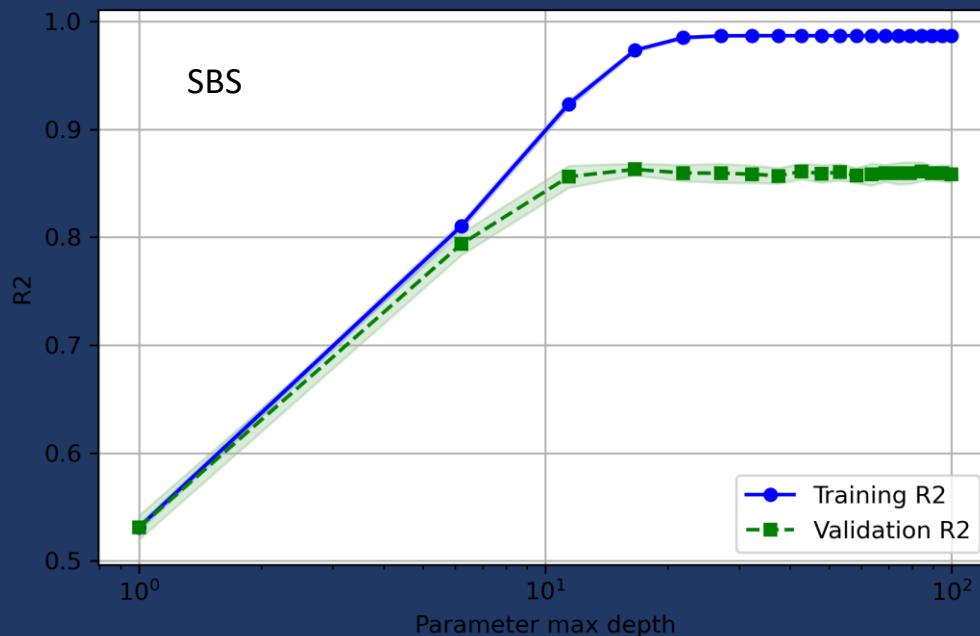
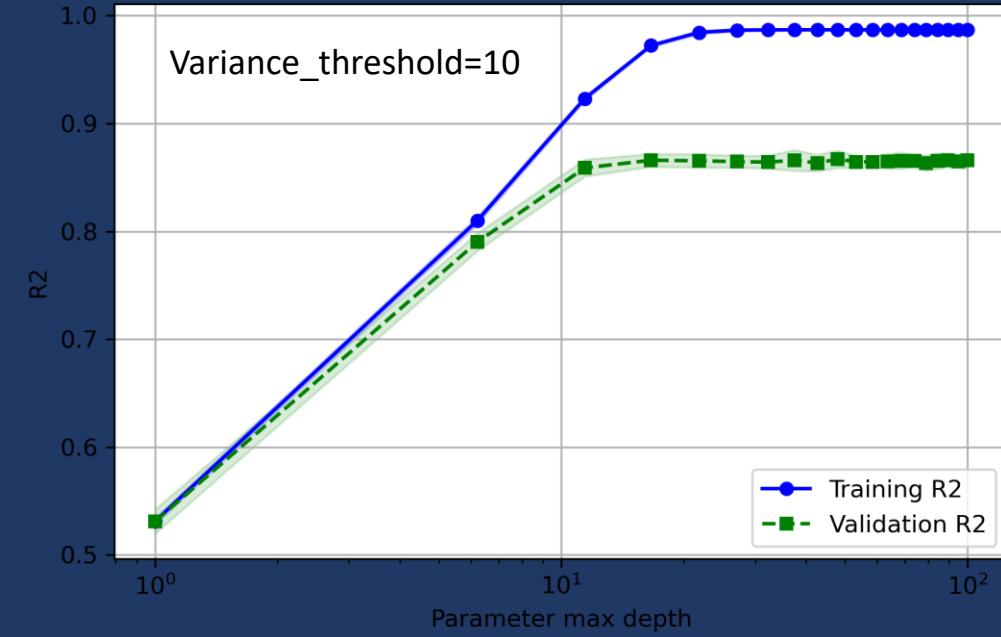
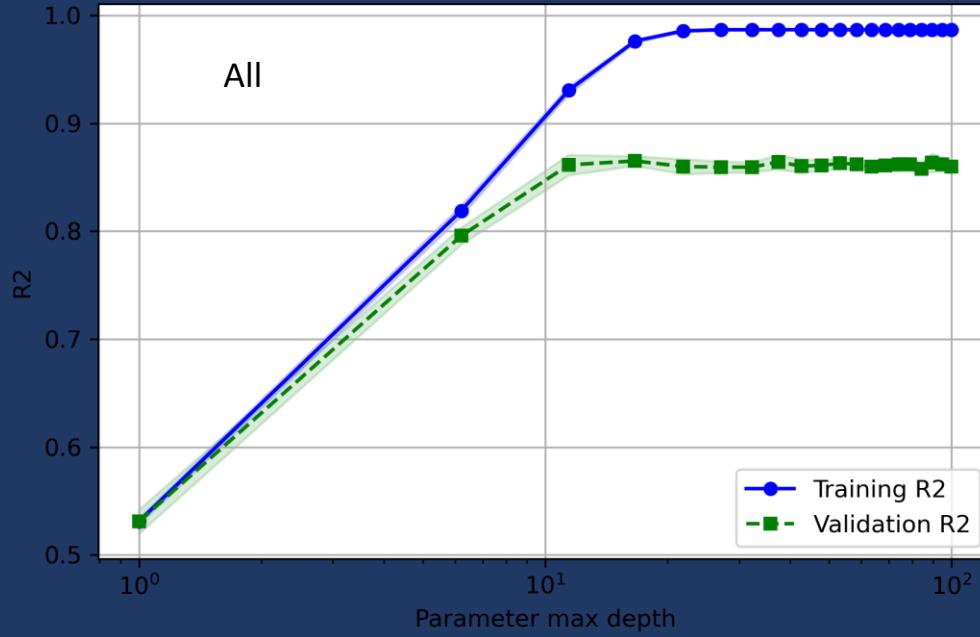
$$I(t) = \frac{1}{N_t} \sum_{i \in D_t}^c (y^{(i)} - \hat{y}_t)^2$$

PROS	CONS
<ul style="list-style-type: none"><li>• Understanding</li><li>• No standardization needed</li><li>• Low Computational cost</li></ul>	<ul style="list-style-type: none"><li>• Overfitting</li><li>• High noise sensibility</li><li>• Hyperparameters</li></ul>

## PARAMETERS

- max\_depth: The maximum depth of the tree
- criterion: The function to measure the quality of a split
- splitter: The strategy used to choose the split at each node
- min\_samples\_split: The minimum number of samples required to split an internal node:
- min\_samples\_leaf: The minimum number of samples required to be at a leaf node.

BEST PARAMETERS: max\_depth=12,, criterion= 'squared\_error', splitter= ' best'.



# Decision Tree

Bagging

Random  
Forest

Boosting

Extreme Gradient  
Boosting

# Random Forest Regressor

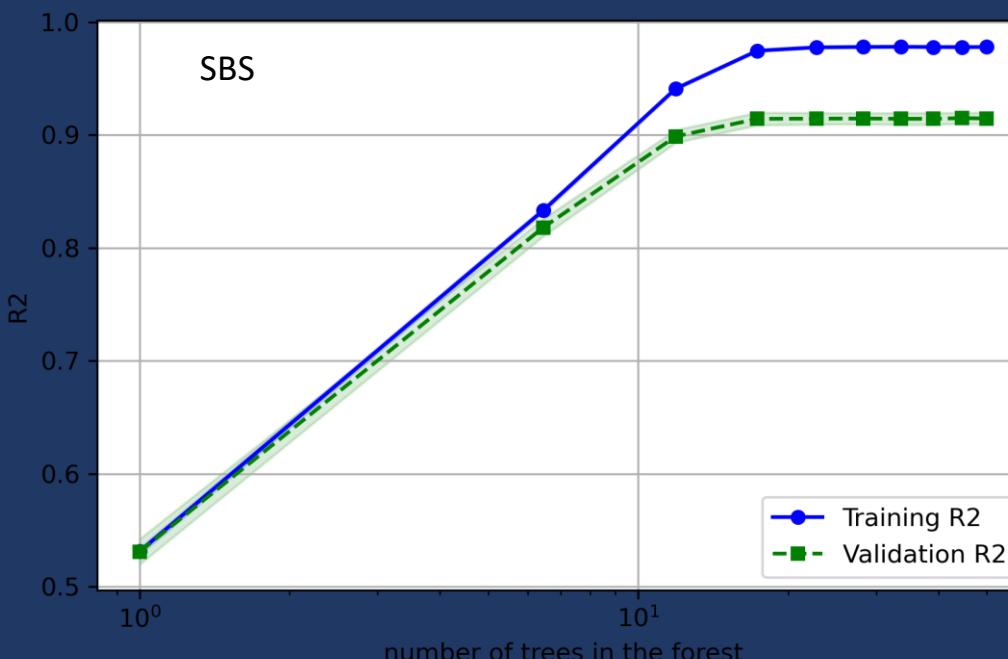
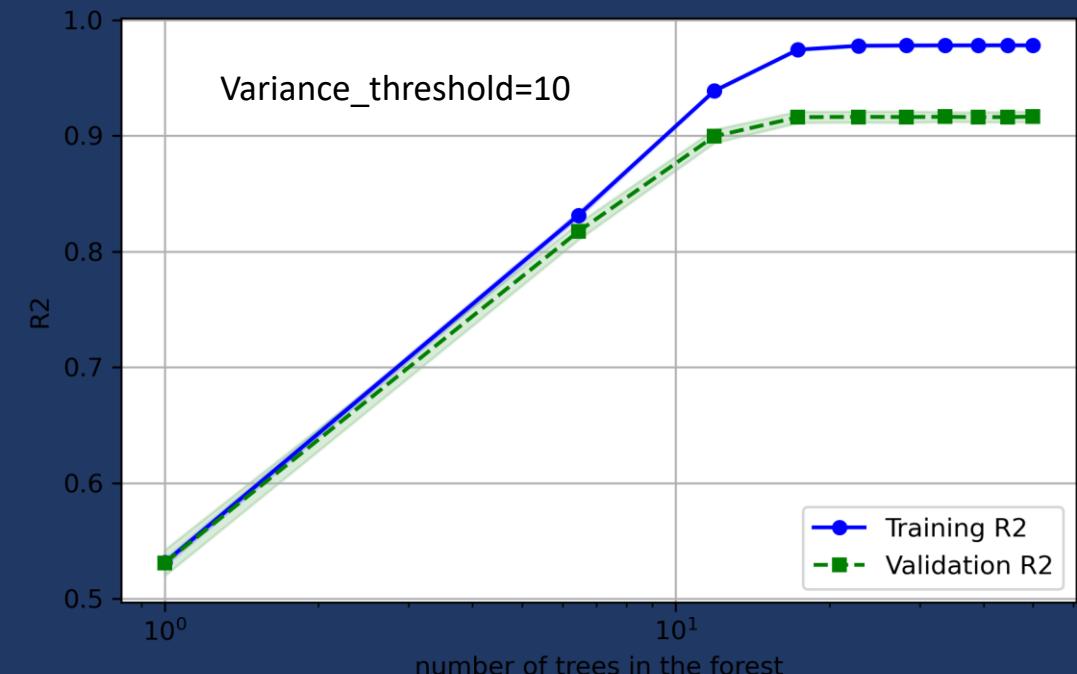
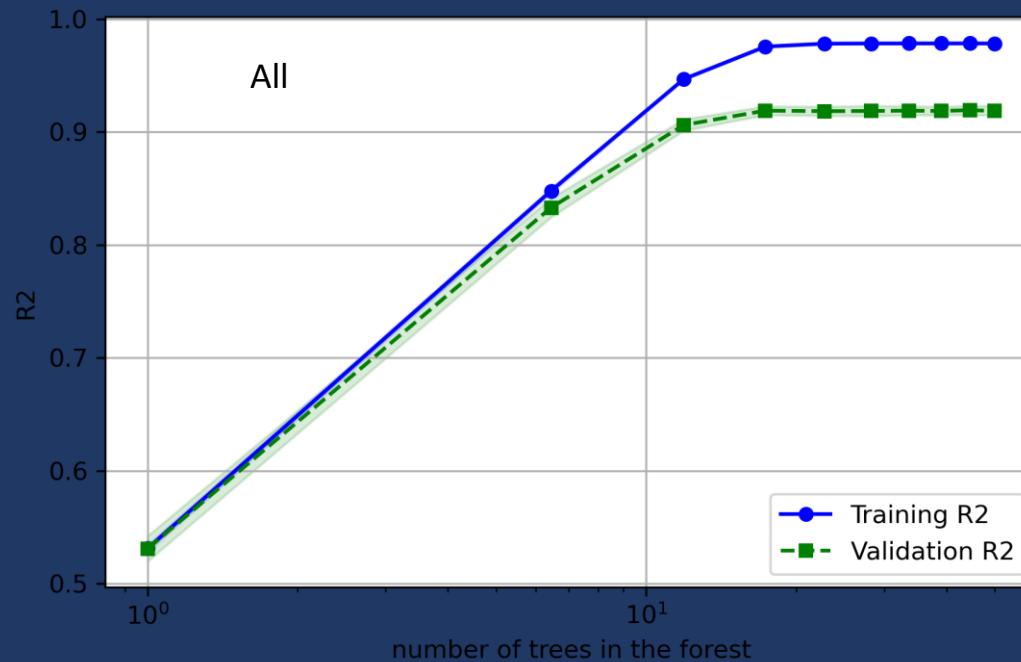
- Ensamble Method (Bagging)
- Decision Tree “in parallel”

PROS	CONS
<ul style="list-style-type: none"><li>• Invariant scale</li><li>• Best generalization performance</li><li>• Simple tuning</li></ul>	Overfitting

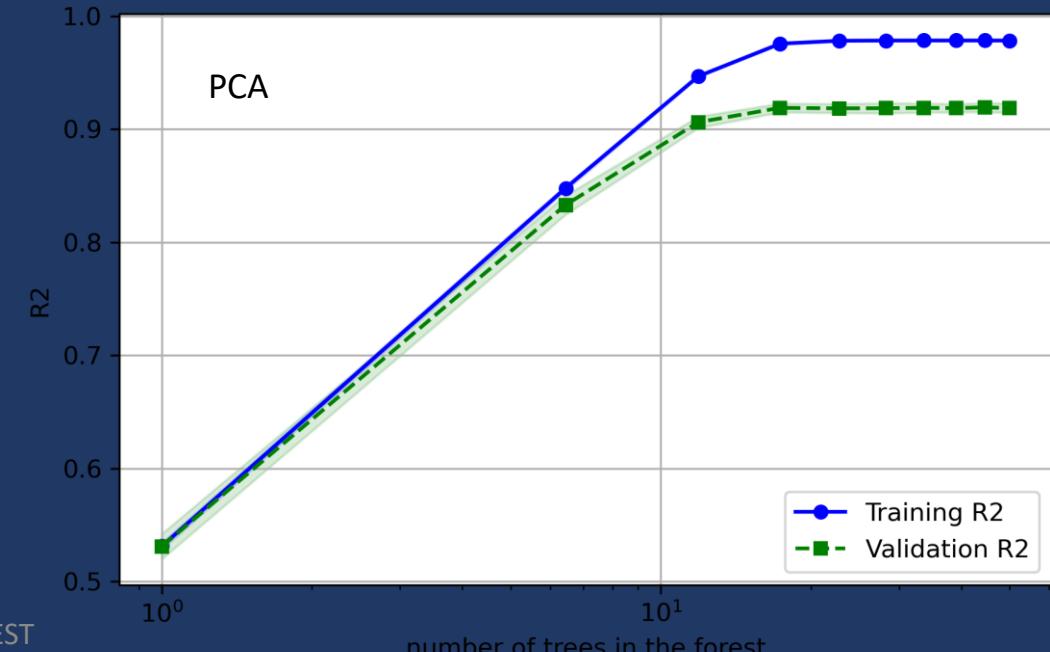
## PARAMETERS

- **n\_estimators:** The number of trees in the forest
- **max\_depth:** The maximum depth of the tree
- **criterion:** The function to measure the quality of a split
- **min\_samples\_split:** The minimum number of samples required to split an internal node:
- **min\_samples\_leaf:** he minimum number of samples required to be at a leaf node.

BEST PARAMETERS: n\_estimators=100, max\_depth = 15, criterion= 'squared\_error'.



RANDOM FOREST



# Extreme Gradient Boosting Regressor

- Ensemble Method (Boosting)
- Gradient boosting involves three elements:
  1. A loss function (MSE) to be optimized.
  2. A weak learner( tree) to make predictions.
  3. An additive model to add weak learners to minimize the loss function.

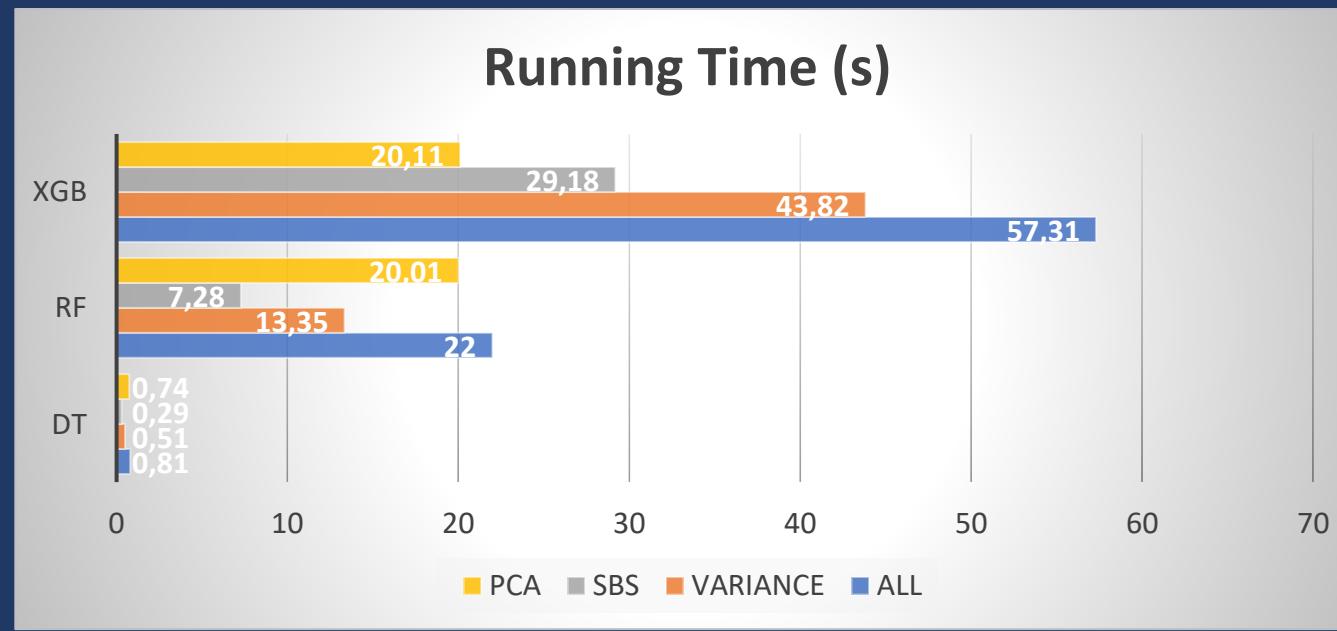
PROS	CONS
<ul style="list-style-type: none"><li>• Reduce Variance and Bias with respect to DT</li></ul>	<ul style="list-style-type: none"><li>• Computationally hard</li></ul>

## PARAMETERS

- **n\_estimators:** Number of gradient boosted trees
- **eta:** Step size shrinkage used in update to prevents overfitting
- **max\_depth:** Maximum depth of a tree
- **subsample:** Subsample ratio of the training instances

BEST PARAMETERS: n\_estimators=350, eta = 0,02, max\_depth = 15, subsample= 0,5. (Grid Search)

	METRICS	ALL	VARIANCE	SBS	PCA
DT	R2 TRAIN	0,94	0,94	0,94	0,94
	R2 TEST	0,86	0,86	0,86	0,86
RF	R2 TRAIN	0,97	0,97	0,97	0,97
	R2 TEST	0,92	0,91	0,91	0,92
XGB	R2 TRAIN	0,98	0,98	0,98	0,97
	R2 TEST	0,92	0,92	0,93	0,92



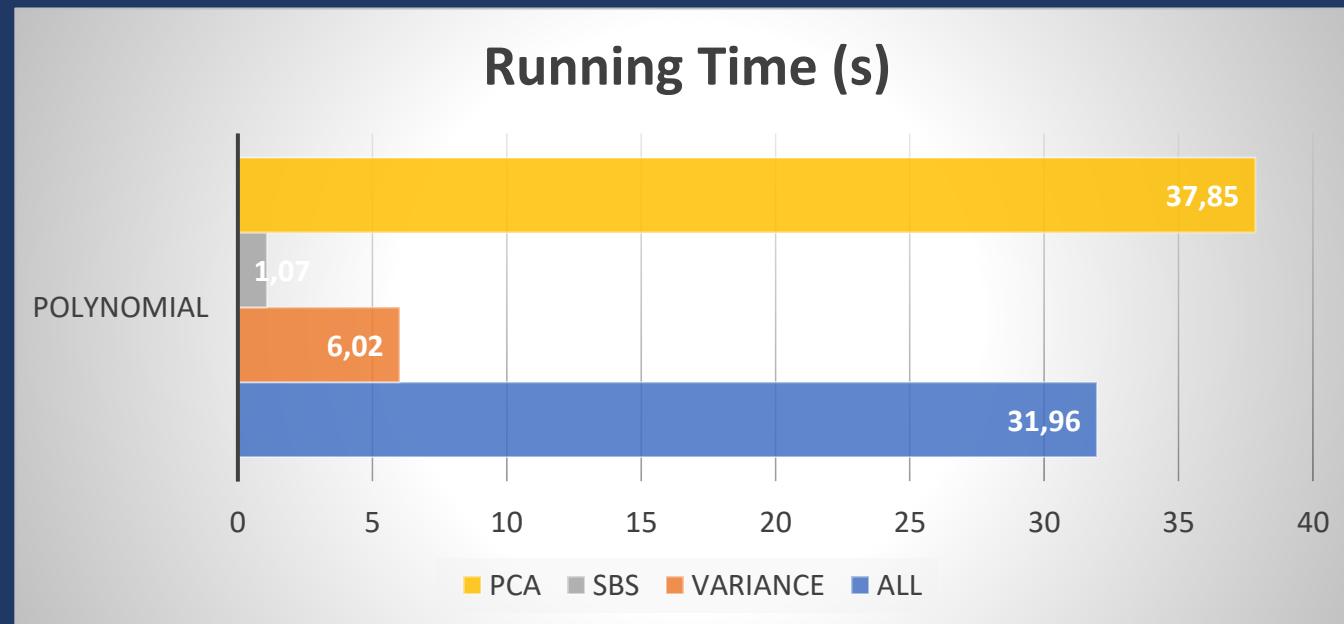
# Multiple Polynomial Regression

- This algorithm allows us to describe the data with a polynomial of degree n :

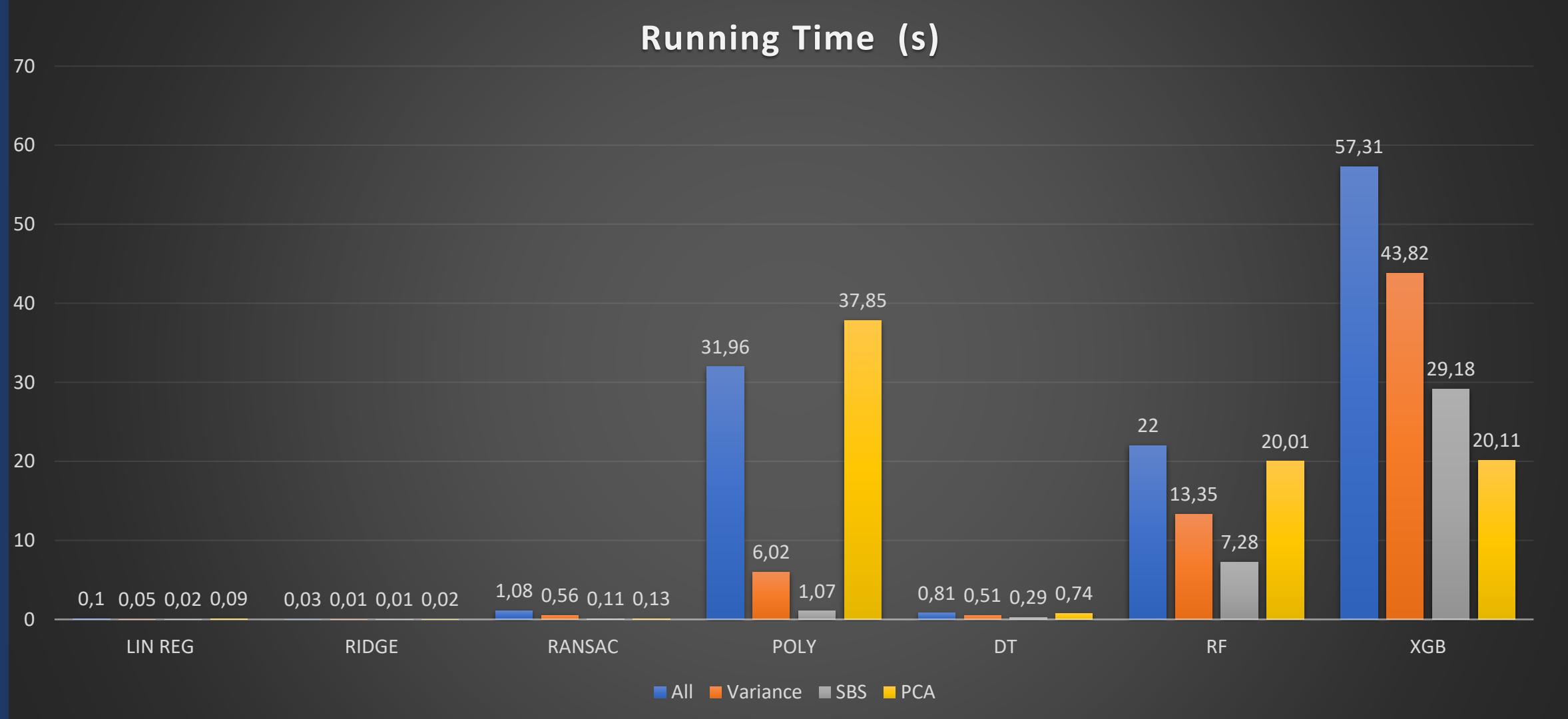
$$y = w_0 + w_1x + w_2x^2 + \cdots + w_nx^n$$

- Together with this model the algorithm was quite efficient.
- On Superconductors dataset, the best degree for the polynomial fitting curve is **n=2**.

METRICS	ALL	VARIANCE	SBS	PCA
R2 TRAIN	0,91	0,87	0,84	0,91
R2 TEST	0,69	0,82	0,81	0,69



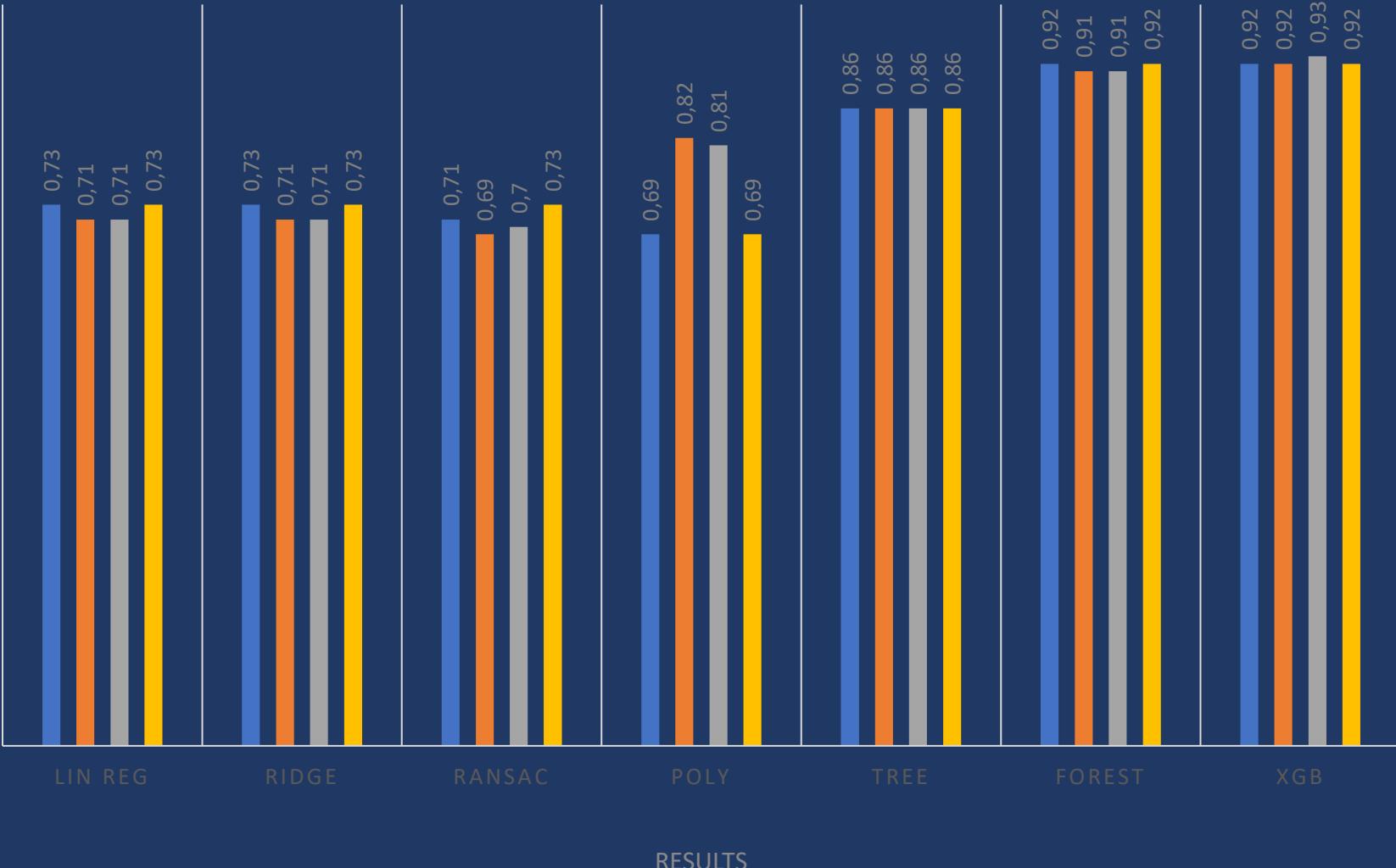
# Results



# Results

## TEST SCORE (R2)

■ ALL ■ VARIANCE ■ SBS ■ PCA



# Conclusions

- We have managed to achieve resourceful results with most models, but the best in terms of performance is the Extreme Gradient Boost.
- The model that represents the best compromise between algorithm execution time and results is the Random Forest Regressor on the features obtained through SBS .
- It is important to note that with this model **it is not possible to predict whether a material is superconductive or not**. However, if we have a superconductor material, we can estimate its critical temperature by knowing some chemical-physical characteristics.
- Taking RMSE as error estimate, we can say that we are able to predict the critical temperature of a superconductor with an **error of about 9,4 Kelvin**.