

Machine Learning Meets (Quantum) Chemistry: Introduction and Short Overview

K. Briling, P. van Gerwen, A. Fabrizio, B. Meyer, R. Fabregat

École Polytechnique Fédérale de Lausanne
Laboratory for Computational Molecular Design (LCMD)
Prof. Clémence Corminboeuf

ksenia.briling@epfl.ch

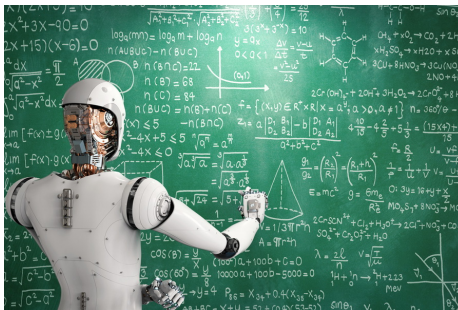
12 May 2025



Introduction



Can a computer learn (quantum) chemistry?





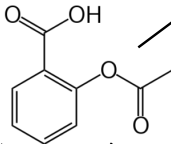
- For example, learn to predict the (quantum) chemical properties from molecular structural formulas

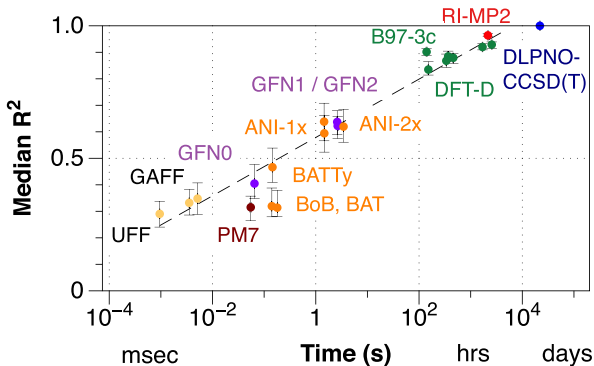
Pharmacological
properties ?

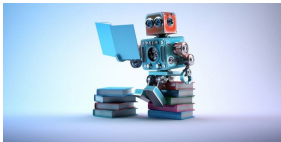
How does it react
with another
compound?

How are its spectra
(Infrared, NMR,...)

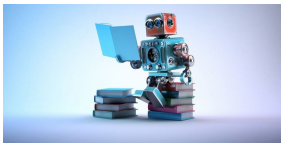
Is its solid
form stable?



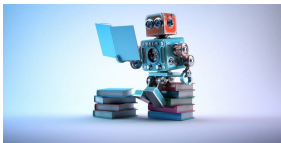




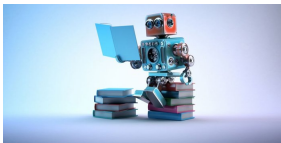
- A computer can learn, just like first experimental chemists have typically learned chemistry...



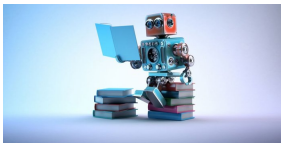
- A computer can learn, just like first experimental chemists have typically learned chemistry...
- **Empirically**: Based on their experience as opposed to *a priori* theoretical knowledge



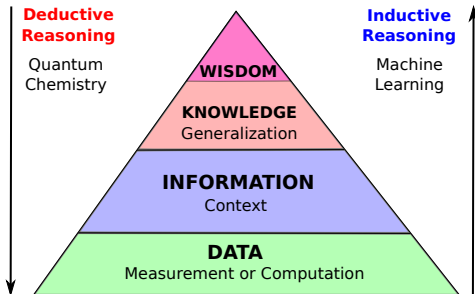
- A computer can learn, just like first experimental chemists have typically learned chemistry...
- **Empirically**: Based on their experience as opposed to *a priori* theoretical knowledge
- Starting from a dataset of theoretical or experimental data with molecular structures and the corresponding observable properties



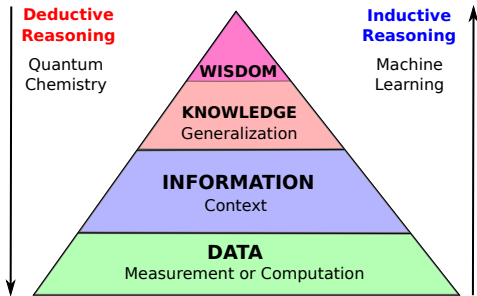
- A computer can learn, just like first experimental chemists have typically learned chemistry...
- **Empirically**: Based on their experience as opposed to *a priori* theoretical knowledge
- Starting from a dataset of theoretical or experimental data with molecular structures and the corresponding observable properties
- The computer can find relationships between molecular structures and the properties



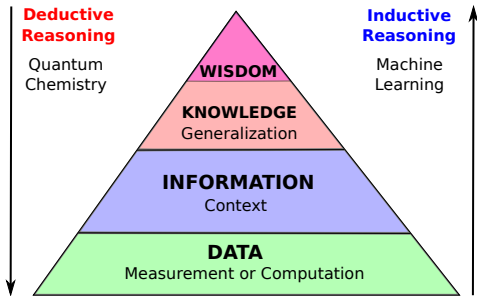
- A computer can learn, just like first experimental chemists have typically learned chemistry...
- **Empirically**: Based on their experience as opposed to *a priori* theoretical knowledge
- Starting from a dataset of theoretical or experimental data with molecular structures and the corresponding observable properties
- The computer can find relationships between molecular structures and the properties
- It learns! And can apply the knowledge to new situations!



- ML is traditionally used in quantitative structure property (or activity) relationships (QSPR/QSAR)



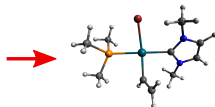
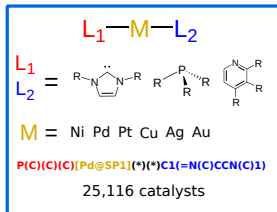
- ML is traditionally used in quantitative structure property (or activity) relationships (QSPR/QSAR)
- Can we unify the two worlds to extend the current quantum chemistry toolbox by “quantum machine learning” (QML)?



- ML is traditionally used in quantitative structure property (or activity) relationships (QSPR/QSAR)
- Can we unify the two worlds to extend the current quantum chemistry toolbox by “quantum machine learning” (QML)?
- QML: classical machine learning applied to quantum-chemical properties

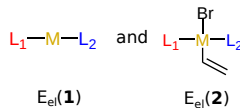


Catalysts Library

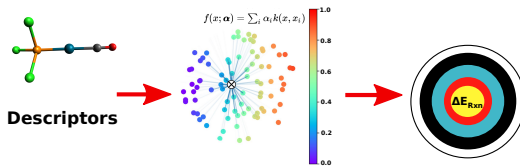


3D Geometry

ab initio Geometries
and Energies for



Machine Learning Models



$$\Delta E_{\text{Rxn}} = E_{\text{el}}(2) - E_{\text{el}}(1)$$

Machine learning methods: a short overview



- Supervised learning (*Classification, Regression*)

Find relations between your data (molecules) and a target variable (chemical properties) that you want to be able to predict



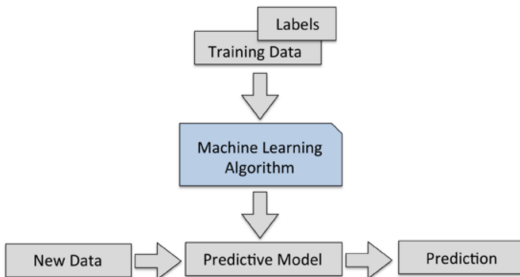
- Supervised learning (*Classification, Regression*)
Find relations between your data (molecules) and a target variable (chemical properties) that you want to be able to predict
- Unsupervised learning (*Clustering, Dimensionality Reduction*)
Searching for indirect hidden structures, clusters, patterns or features in the raw (molecular) data



- Supervised learning (*Classification, Regression*)
Find relations between your data (molecules) and a target variable (chemical properties) that you want to be able to predict
- Unsupervised learning (*Clustering, Dimensionality Reduction*)
Searching for indirect hidden structures, clusters, patterns or features in the raw (molecular) data
- Reinforcement learning
Solving interactive problems with an environment (e.g., a chess engine)
Close to supervised learning but reward (feedback) instead of labels

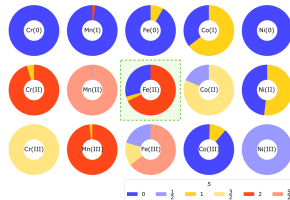
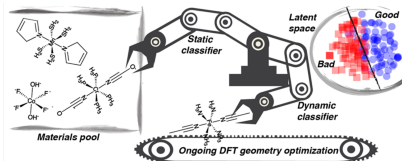
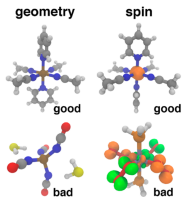


- In supervised learning the model is constructed from *training* molecules with known properties that allows us to make predictions about unseen molecules
- *Classification* task: assign discrete class labels
- *Regression* task: the chemical property is a continuous value





- Goal: predict the categorical labels of new molecules based on past observations
- Applications: electronic structure computations outcomes (good/bad), ground state spins, ...



C. Duan et al., "Learning from failure: Predicting electronic structure calculation outcomes with machine learning models", *J. Chem. Theory Comput.* **15**, 2331–2345 (2019); Y. Cho et al., "Automated prediction of ground state spin for transition metal complexes", *Digital Discovery* **3**, 1638–1647 (2024)



- Task: prediction of continuous properties
- Applications: heat capacity at room temperature, HOMO–LUMO gap, receptor–ligand binding, stability of molecular conformers, ...

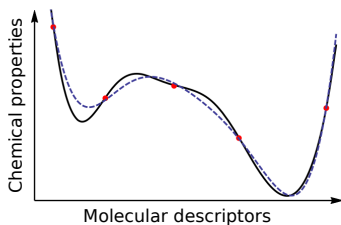


- Task: prediction of continuous properties
- Applications: heat capacity at room temperature, HOMO–LUMO gap, receptor–ligand binding, stability of molecular conformers, ...

Can be applied successfully if:

- Cause and effect relationship connecting system to property
- Query scenario is interpolative in nature
- Sufficient training data available

$$\mathbf{Y} = \mathbf{f}(\text{descriptors})$$

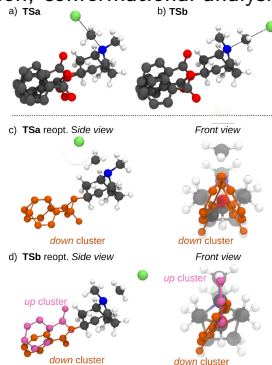
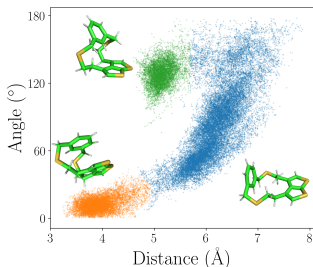


Unsupervised ML:

Finding subgroups with clustering

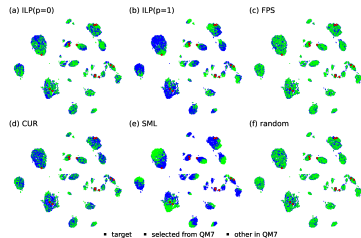
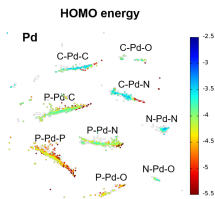
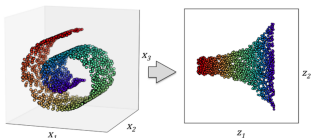


- Exploratory data analysis technique that allows to organize information into meaningful subgroups (clusters) without any prior knowledge
- Each cluster defines a group of objects that share a certain degree of similarity but are more dissimilar to objects in other clusters
- Applications: representative sample, subsets selection, conformational analysis





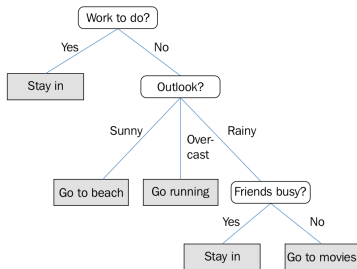
- Clustering identifies agglomeration of data but do not offer an overall picture of the relations between different structures
- Dimensionality reduction compresses the data onto a smaller dimensional subspace while retaining most of the relevant information
- Data visualization: a high-dimensional feature set can be projected onto 1-, 2-, or 3D spaces
- Applications: representation of chemical space, pattern recognition



B. Sawatlon et al., "Data mining the C–C cross-coupling genome", *ChemCatChem* **11**, 4096–4107 (2019); M. Haeberle et al., "Integer linear programming for unsupervised training set selection in molecular machine learning", *Mach. Learn.: Sci. Technol.* **6**, 025030 (2025)

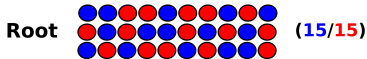


- Decision Trees (DTs) are a **non-parametric** supervised learning method used for classification and regression
- Learns simple decision rules inferred from the data features
- DT split a set of objects into subsets (usually 2 in binary trees) that are purer in composition
- Rules are selected based on how well splits can differentiate





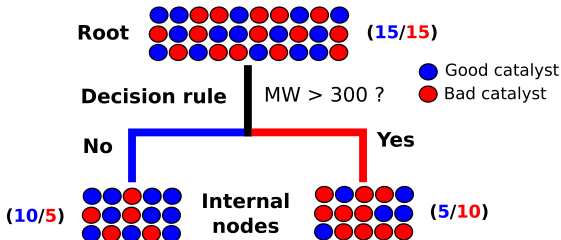
- Starting with a root composed by a set of 30 catalysts and knowing the ones that are “good” or “bad” for a given reaction
- Each catalyst is described by several simple molecular descriptors (e.g., molecular weight, shape index, molecular volume, # of C atoms, # of rotatable bonds. . .



- Good catalyst
- Bad catalyst

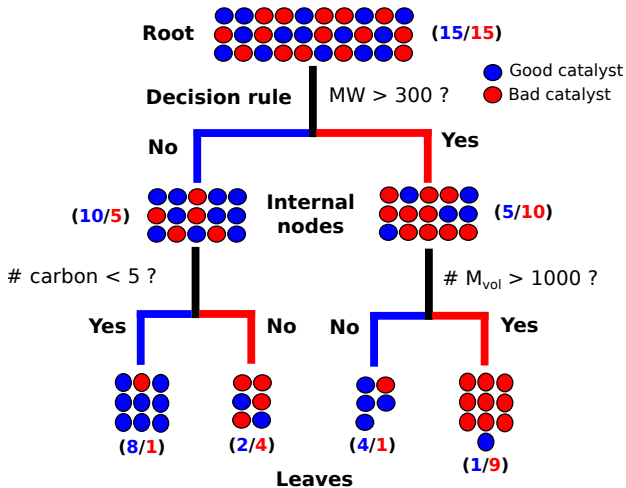


- DT algorithm searches a rule that yields the highest decrease in impurity
- Once a rule is selected and splits a node into two





- The same logic is applied to each “child” node
- Stops when no further gain / stopping rules met





Why classification and regression trees (CART) is a successful tool?



Why classification and regression trees (CART) is a successful tool?

- Universally applicable to classification and regression problems with no assumptions on the data structure
- The picture of the tree structure gives valuable insights into which variables are important
- Terminal nodes give a natural clustering of the data into homogenous groups
- Can handle large data sets: $O(DM \log M)$
(M is # of molecules, D is # of descriptors)



Why classification and regression trees (CART) is a successful tool?

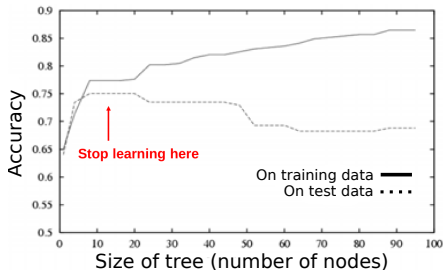
- Universally applicable to classification and regression problems with no assumptions on the data structure
- The picture of the tree structure gives valuable insights into which variables are important
- Terminal nodes give a natural clustering of the data into homogenous groups
- Can handle large data sets: $O(DM \log M)$
(M is # of molecules, D is # of descriptors)

But...

- Decision-tree learners can create over-complex trees that do not generalize the data well
- This is called **overfitting**

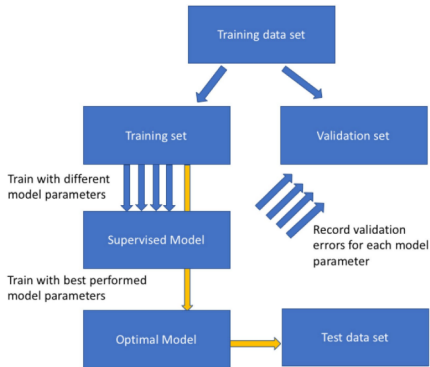
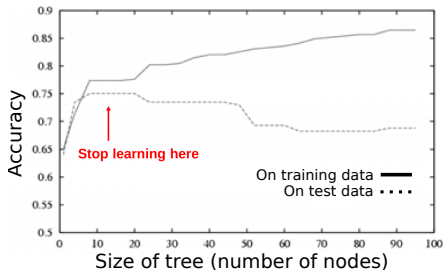


- Decision trees can overfit data
- So, it is necessary to use a **validation** set in order to prune the tree at an optimal size



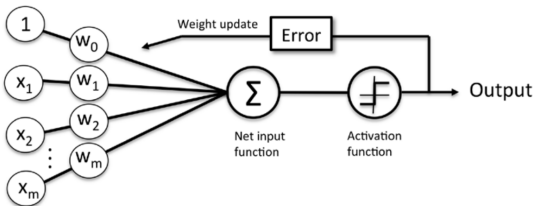


- Decision trees can overfit data
- So, it is necessary to use a **validation** set in order to prune the tree at an optimal size



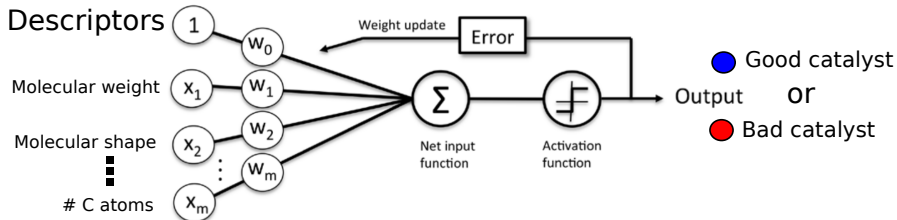


- Artificial neurons are inspired from biological neurons
- An artificial neuron has one or many inputs, each associated to a weight
- If the weighted sum of inputs is lower than a threshold, the neuron remains inactive
- If the weighted sum bypass the threshold, the neuron is activated and produce an output signal
- During the learning phase, this output is used to calculate the error of the prediction and update the weights



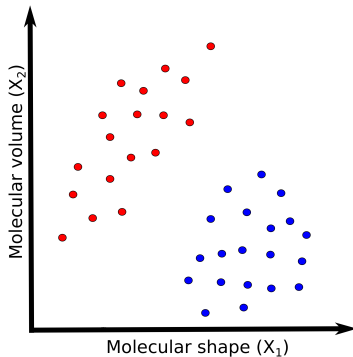


Starting with the same set of 30 catalysts and knowing the ones that are “good” or “bad” for a given reaction



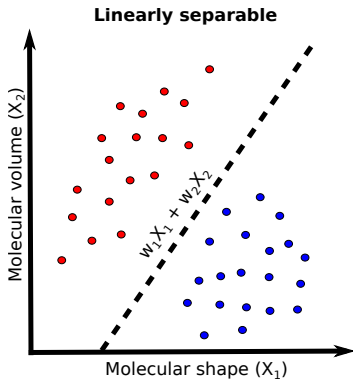


During the learning phase, the outputs are used to update the weights



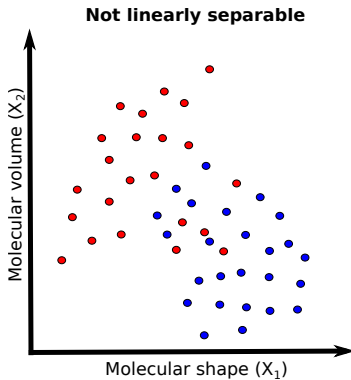


If the two classes are linearly separable, the perceptron will converge to a solution



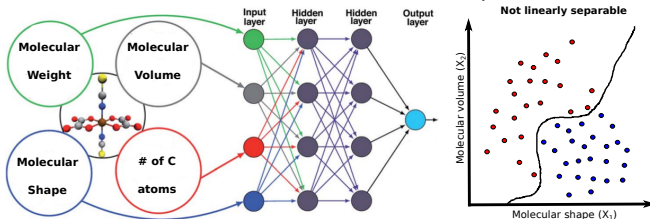


However, if the two classes are not linearly separable, either we accept the error, either we need more than a single neuron

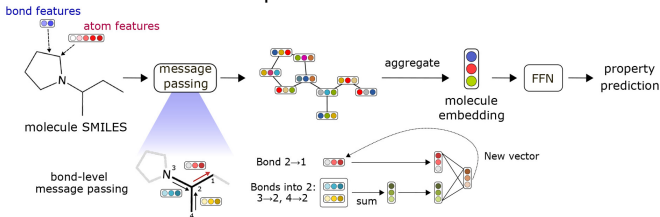




- Learns non-linear models for classification and prediction



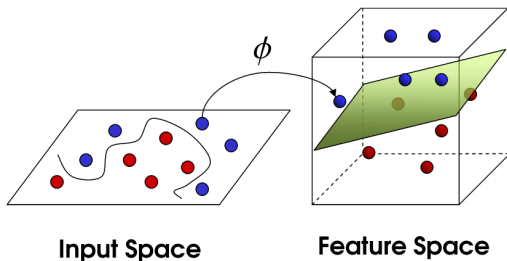
- Requires tuning a large number of hyperparameters
- Can learn new data descriptors itself





Another way to solve non-linear problem

- Any non-linear problem can be mapped (ϕ) into a higher-dimensional feature space where it becomes linearly separable
- Computationally very expensive
- The step $\phi(\mathbf{x}) \cdot \phi(\mathbf{x}')$ is replaced by a kernel function $K(\mathbf{x}, \mathbf{x}')$

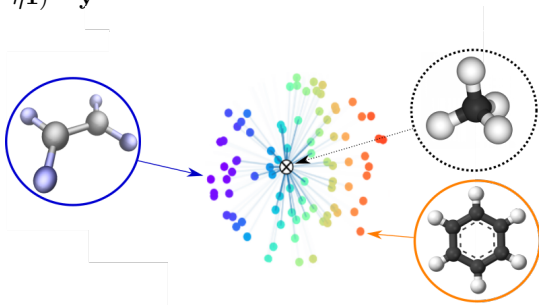




- Roughly speaking, kernels can be interpreted as similarity measures between pair of molecules
- The most widely used kernels are Gaussian $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma^2(\mathbf{x} - \mathbf{x}')^2)$ and Laplacian $K(\mathbf{x}, \mathbf{x}') = \exp\left(-\gamma \sum_d^D |x_d - x'_d|\right)$



- Roughly speaking, kernels can be interpreted as similarity measures between pair of molecules
- The most widely used kernels are Gaussian $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma^2(\mathbf{x} - \mathbf{x}')^2)$ and Laplacian $K(\mathbf{x}, \mathbf{x}') = \exp\left(-\gamma \sum_d^D |x_d - x'_d|\right)$
- The prediction for y_t is a weighted sum of kernel functions centered on each training point: $y_t = \sum_i^N K(\mathbf{x}_t, \mathbf{x}_i) w_i^*$
- The regression coefficients represent the contribution of each training point to the target value and minimize the quadratic loss function on the training set: $\mathbf{w}^* = (\mathbf{K} + \eta \mathbf{1})^{-1} \mathbf{y}$



Molecular representations



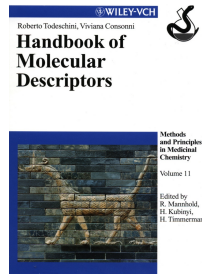
“... the final result of a logic and mathematical procedure which transforms chemical information encoded within a symbolic representation of a molecule into a useful number”



“... the final result of a logic and mathematical procedure which transforms chemical information encoded within a symbolic representation of a molecule into a useful number”

Vast majority:

- constitutional (# of atoms, MW,...)
- graph-based (bond-, distance-, adjacency matrices,...)
- mainly used in bio or chemo-informatics for QSAR applications





- Focus on QML models: physics-based, systematic, universal



- Focus on QML models: physics-based, systematic, universal
- To better distinguish QML from QSAR, we prefer “representation” to “descriptor”



- Focus on QML models: physics-based, systematic, universal
- To better distinguish QML from QSAR, we prefer “representation” to “descriptor”
- The ground state of a chemical system is defined by its Hamiltonian



- Focus on QML models: physics-based, systematic, universal
- To better distinguish QML from QSAR, we prefer “representation” to “descriptor”
- The ground state of a chemical system is defined by its Hamiltonian
- Hamiltonian depends on elemental composition and geometry (and number of electrons)

$$\hat{H}(\mathbf{Z}_i, \mathbf{R}_i) \xrightarrow{\text{QM}} E$$



- Focus on QML models: physics-based, systematic, universal
- To better distinguish QML from QSAR, we prefer “representation” to “descriptor”
- The ground state of a chemical system is defined by its Hamiltonian
- Hamiltonian depends on elemental composition and geometry (and number of electrons)

$$\hat{H}(\mathbf{Z}_i, \mathbf{R}_i) \xrightarrow{\text{QM}} E$$

- Representation in ML plays the role of Hamiltonian/wavefunction in QM



- Focus on QML models: physics-based, systematic, universal
- To better distinguish QML from QSAR, we prefer “representation” to “descriptor”
- The ground state of a chemical system is defined by its Hamiltonian
- Hamiltonian depends on elemental composition and geometry (and number of electrons)

$$\hat{H}(\mathbf{Z}_i, \mathbf{R}_i) \xrightarrow{\text{QM}} E$$

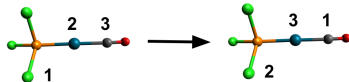
- Representation in ML plays the role of Hamiltonian/wavefunction in QM
- It should be a vector which encodes composition and geometry (and charge/spin) of a molecule

$$\text{representation}(\mathbf{Z}_i, \mathbf{R}_i) \xrightarrow{\text{ML}} E$$



Criteria:

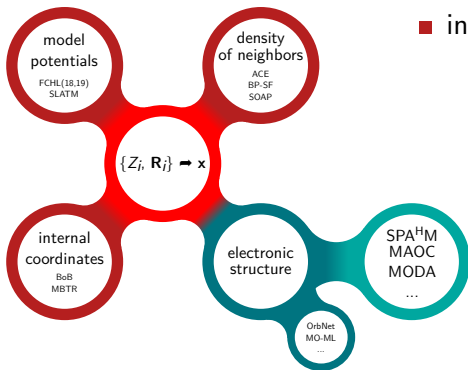
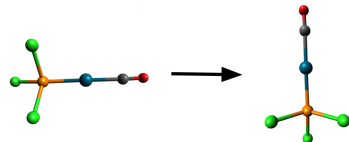
- injectivity
- continuity and differentiability
- invariance/equivariance w.r.t.:
 - permutations



- translations



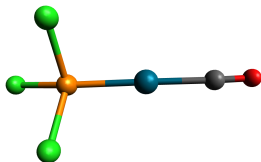
- rotations





- Square atom-by-atom matrix
- Off-diagonal elements correspond to the Coulomb repulsion between nuclei
- Diagonal elements remind the electronic energy of a H-like atom

$$M_{ij} = \begin{cases} \frac{Z_i Z_j}{|\mathbf{R}_i - \mathbf{R}_j|} & \text{for } i \neq j \\ 0.5 Z_i^{2.4} & \text{for } i = j \end{cases}$$

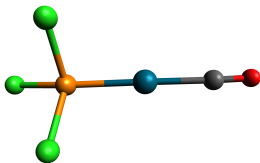


	P	Cl	Cl	Cl	Pd	C	O
P	332	124	122	122	290	20	22
Cl	124	449	94	94	179	16	18
Cl	122	94	449	82	232	20	22
Cl	122	94	82	449	237	20	22
Pd	290	179	232	237	4893	136	116
C	20	16	20	20	136	37	43
O	22	18	22	22	116	43	74



- Square atom-by-atom matrix
- Off-diagonal elements correspond to the Coulomb repulsion between nuclei
- Diagonal elements remind the electronic energy of a H-like atom

$$M_{ij} = \begin{cases} \frac{Z_i Z_j}{|\mathbf{R}_i - \mathbf{R}_j|} & \text{for } i \neq j \\ 0.5 Z_i^{2.4} & \text{for } i = j \end{cases}$$

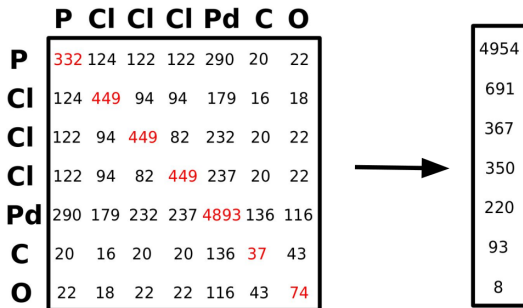


	P	Cl	Cl	Cl	Pd	C	O
P	332	124	122	122	290	20	22
Cl	124	449	94	94	179	16	18
Cl	122	94	449	82	232	20	22
Cl	122	94	82	449	237	20	22
Pd	290	179	232	237	4893	136	116
C	20	16	20	20	136	37	43
O	22	18	22	22	116	43	74

- No well-defined ordering of the atoms in the matrix



- Represent the molecule as a vector of sorted eigenvalues of the CM
- Invariant with respect to atom permutations
- Loss of information





- Permute the matrix in order to sort the rows and the columns by their norm
- Unique CM representation
- More information than eigenspectrum

	P	Cl	Cl	Cl	Pd	C	O
P	332	124	122	122	290	20	22
Cl	124	449	94	449	179	16	18
Cl	122	94	449	82	232	20	22
Cl	122	94	82	449	237	20	22
Pd	290	179	232	237	4893	136	116
C	20	16	20	20	136	37	43
O	22	18	22	22	116	43	74



	Pd	Cl	Cl	Cl	P	C	O
Pd	4893	237	232	179	290	20	22
Cl	237	449	82	94	122	16	18
Cl	232	82	449	94	122	20	22
Cl	179	94	94	449	124	20	22
P	290	122	122	124	332	136	116
C	20	16	20	20	136	37	43
O	22	18	22	22	116	43	74



- Construct several CMs based on a random ordering of the atoms, adding a random noise ϵ to the row norms $\|C\|$ and determine the permutation that minimizes $\|C\| + \epsilon$
- Approximate sampling of all possible valid CMs given a specific molecule
- Increased the computational costs

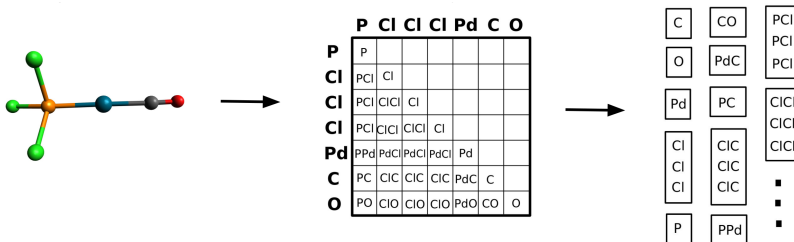
	P	Cl	Cl	Cl	Pd	C	O
P	332	124	122	122	290	20	22
Cl	124	449	94	94	179	16	18
Cl	122	94	449	82	232	20	22
Cl	122	94	82	449	237	20	22
Pd	290	179	232	237	4893	136	116
C	20	16	20	20	136	37	43
O	22	18	22	22	116	43	74



	P	Cl	Cl	Cl	Pd	C	O					
P	332	124	122	122	290	20	22					
Cl	124											
Cl	122	P	332	124	122	290	20	22				
Cl	122	Cl	124	4								
Pd	290	Cl		Pd	4893	237	232	179	290	20	22	
C	20	Cl	122		Cl	237	449	82	94	122	16	18
O	22	Pd	290	1	Cl	232	82	449	94	122	20	22
		C	20		Cl	179	94	94	449	124	20	22
		O	22		P	290	122	122	124	332	136	116
					C	20	16	20	20	136	37	43
					O	22	18	22	22	116	43	74



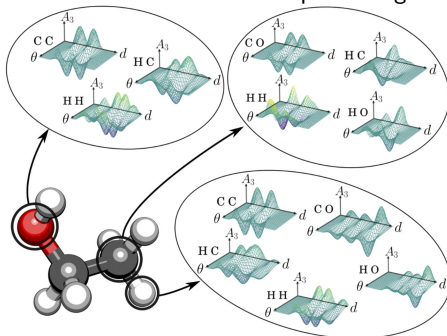
- Inspired by NLP: bag-of-words descriptor encodes the frequency of occurrence of words in text
- Each bag corresponds to a specific type of atomic pair
- For example, all P-Cl pairs in the molecule are grouped into the bag labeled as PCI
- Crucial higher-order information (angles and dihedrals) missing



SLATM: the Spectrum of London and Axilrod–Teller–Muto potential



- Represents an atom i by accounting for all possible interactions between it and its neighboring atoms through many-body potential terms multiplied by a Gaussian distribution G and put in bags



SLATM: the Spectrum of London and Axilrod–Teller–Muto potential

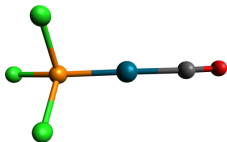


- Represents an atom i by accounting for all possible interactions between it and its neighboring atoms through many-body potential terms multiplied by a Gaussian distribution G and put in bags
- The **one-body** term: simply the nuclear charge ($x_1 = Z_i$)
- The **two-body** term:

$$x_2(\mathbf{r}) = \frac{1}{2} \sum_{j \neq i}^{N_{\text{at}}} G_{\sigma_2}(\mathbf{r} - \mathbf{R}_{ij}) \cdot \frac{Z_i Z_j}{|\mathbf{r}|^6}$$

- The **three-body** term:

$$x_3(\theta) = \frac{1}{3} \sum_{k \neq j \neq i}^{N_{\text{at}}} G_{\sigma_3}(\theta - \theta_{ijk}) \cdot \frac{Z_i Z_j Z_k [1 + \cos \theta \cos \theta_{jki} \cos \theta_{kij}]}{|\mathbf{R}_{ij}|^3 |\mathbf{R}_{jk}|^3 |\mathbf{R}_{ik}|^3}$$



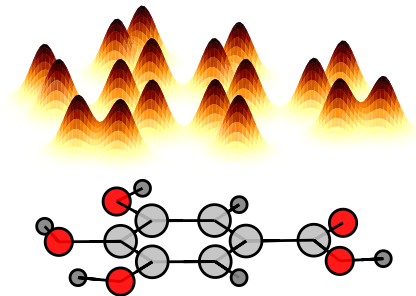
[8], [17], [6], [46], [15], [8, 8], [17, 17], [6, 6], [46, 46], [15, 15], [8, 17], [8, 6], [8, 46], [8, 15], [17, 6], [17, 46], [17, 15], [6, 46], [6, 15], [46, 15], [8, 17, 17], [17, 8, 17], [8, 17, 6], [8, 6, 17], [17, 8, 6], [8, 17, 46], [8, 46, 17], [17, 8, 46], [8, 17, 15], [8, 15, 17], [17, 8, 15], [8, 6, 46], [8, 46, 6], [6, 8, 46], [8, 6, 15], [8, 15, 6], [6, 8, 15], [8, 46, 15], [8, 15, 46], [46, 8, 15], [17, 17, 17], [17, 17, 6], [17, 6, 17], [17, 17, 46], [17, 46, 17], [17, 17, 15], [17, 15, 17], [17, 6, 46], [17, 46, 6], [6, 17, 46], [17, 6, 15], [17, 15, 6], [6, 17, 15], [17, 46, 15], [17, 15, 46], [46, 17, 15], [6, 46, 15], [6, 15, 46], [46, 6, 15]



- Local similarity measure between atoms



- Local similarity measure between atoms
- Each atom i is represented as sum of neighbor nuclei densities smoothened with a Gaussian, $\varrho_i(\mathbf{r}) = \sum_k \exp\left(-\frac{|\mathbf{R}_k - \mathbf{r}|^2}{2\sigma_\varrho^2}\right)$





- Local similarity measure between atoms
- Each atom i is represented as sum of neighbor nuclei densities smoothened with a Gaussian, $\varrho_i(\mathbf{r}) = \sum_k \exp\left(-\frac{|\mathbf{R}_k - \mathbf{r}|^2}{2\sigma_\varrho^2}\right)$
- Similarity is the overlap averaged over rotations

$$\bar{K}_{ij} = \int d\hat{R} \left[\int \varrho_i(\mathbf{r}) \varrho_j(\hat{R}\mathbf{r}) d^3\mathbf{r} \right]^2, \quad K_{ij} = \left[\frac{\bar{K}_{ij}}{\sqrt{\bar{K}_{ii}\bar{K}_{jj}}} \right]^\zeta$$



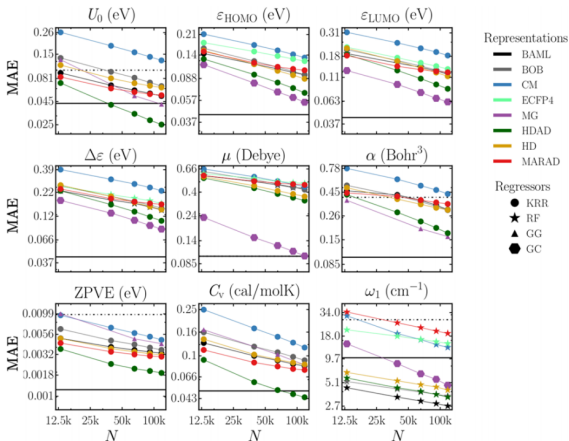
- Local similarity measure between atoms
- Each atom i is represented as sum of neighbor nuclei densities smoothened with a Gaussian, $\varrho_i(\mathbf{r}) = \sum_k \exp\left(-\frac{|\mathbf{R}_k - \mathbf{r}|^2}{2\sigma_\varrho^2}\right)$
- Similarity is the overlap averaged over rotations

$$\bar{K}_{ij} = \int d\hat{R} \left[\int \varrho_i(\mathbf{r}) \varrho_j(\hat{R}\mathbf{r}) d^3\mathbf{r} \right]^2, \quad K_{ij} = \left[\frac{\bar{K}_{ij}}{\sqrt{\bar{K}_{ii}\bar{K}_{jj}}} \right]^\zeta$$

- In practice, ϱ_i is decomposed onto atom-centered basis leading to power spectrum representations \mathbf{p}

$$K_{ij} = \left[\frac{\mathbf{p}_i^\top \mathbf{p}_j}{\sqrt{\mathbf{p}_i^\top \mathbf{p}_i \cdot \mathbf{p}_j^\top \mathbf{p}_j}} \right]^\zeta = [\bar{\mathbf{p}}_i^\top \bar{\mathbf{p}}_j]^\zeta, \quad p_{i,nn'\ell} = \sqrt{\frac{8\pi^2}{2\ell+1}} \sum_m c_{i,n\ell m} c_{i,n'\ell m}^*$$

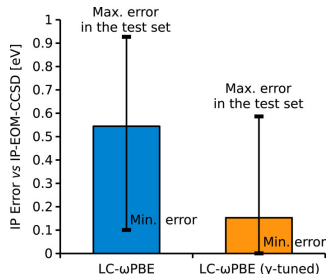
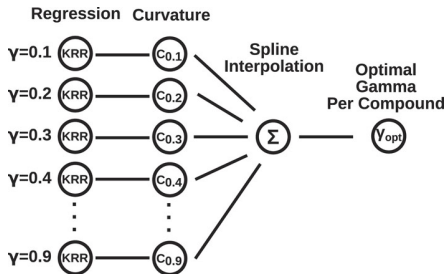
Few examples: ML applied to chemical properties



QM9 database (molecules with up to 9 heavy atoms)

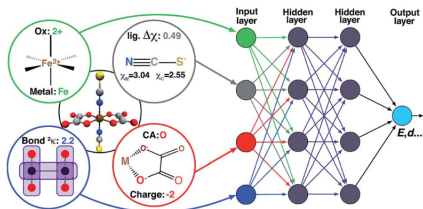


- Average energy curvature is $C_{\text{avg}}^N \int_{N-1}^N \frac{\partial^2 E}{\partial x^2}(x) dx = \varepsilon_{\text{HOMO}}^N - \varepsilon_{\text{LUMO}}^{N-1}$
- Minimization of C is a criterion for the optimal tuning of range-separated hybrid density functionals (γ)
- Pipeline: predict curvature for different γ and find the optimal γ with spline interpolation

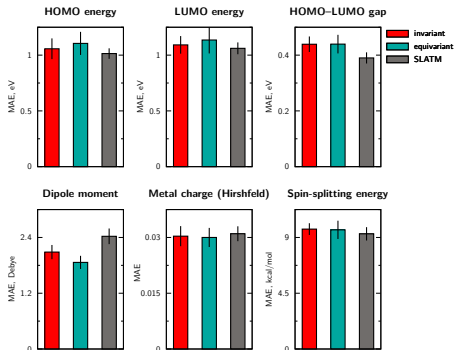




Spin-state ordering, sensitivity to HF exchange, spin-state specific bond lengths

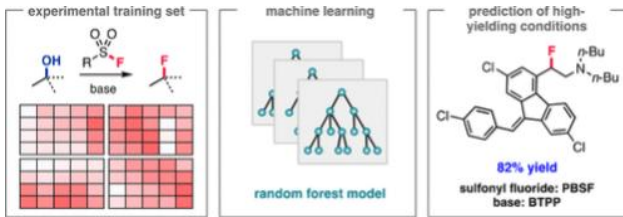


- Graph theory based descriptors: Atomic connectivity, Kier index. Preferred to 3-d structural information (e.g., Coulomb Matrix).
- Complex Based: Metal identity, oxidation state, empirical pairwise Pauling electronegativity, ...
- Atomic descriptors: Charge, mass, ...



- KRR + molecular representations (local, global, structure, electronic...)
- Invariant/equivariant tensor field neural network

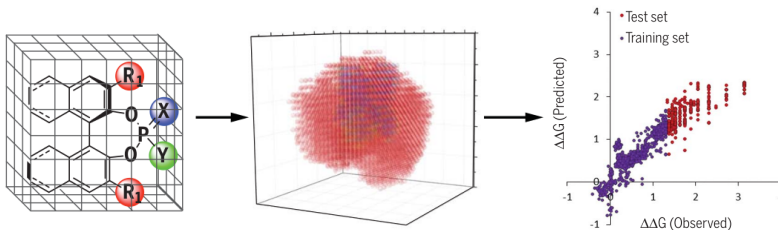
Prediction of the reaction performance in C–N cross-coupling and deoxyfluorination with sulfonyl fluorides from experimental training set constructed by ultra-high-throughput nanoscale experimentation



- Molecular descriptors: Molecular volume, Surface area, Molecular weight, E_{HOMO} , E_{LUMO}
- Atomic descriptors: Electrostatic charge and NMR shift
- Vibrational descriptors: Frequency and intensity

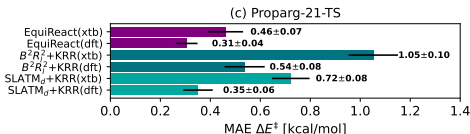
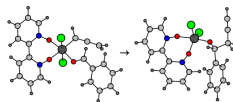
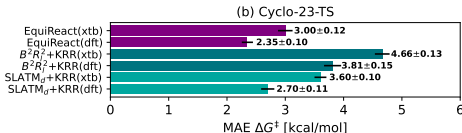
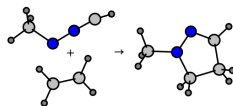
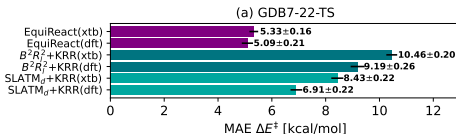
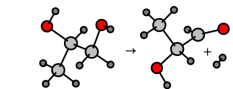


- Exploit machine learning algorithms to accurately predict selective catalysts for the chiral phosphoric acid-catalyzed thiol addition to N-acylimines reactions
- ML models used: Support Vector Machines (SVM) and Neural Network
- Descriptors: average steric occupancy (ASO), computed electrostatic parameters and NBO charges.





- QML for reactions is more challenging
- Involve 2+ components (reactant(s) and product(s))
- Representations design: transition state instead of Hamiltonian

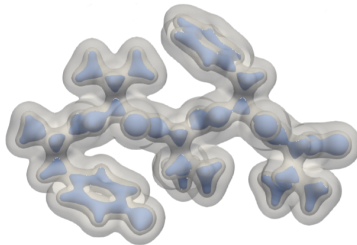




There are many other chemical properties to predict, and not only scalar properties...

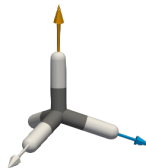
Tensorial properties

Scalar fields



Electron Density

Vectorial fields



Molecular Force Field

Bibliography



Reviews:

- F. Musil et al., “Physics-inspired structural representations for molecules and materials”, *Chem. Rev.* **121**, 9759–9815 (2021)
- J. A. Keith et al., “Combining machine learning and computational chemistry for predictive insights into chemical systems”, *Chem. Rev.* **121**, 9816–9872 (2021)
- B. Huang and O. A. von Lilienfeld, “Ab initio machine learning in chemical compound space”, *Chem. Rev.* **121**, 10001–10036 (2021)
- O. T. Unke et al., “Machine learning force fields”, *Chem. Rev.* **121**, 10142–10186 (2021)

ML books:

- C. E. Rasmussen and C. K. I. Williams, *Gaussian processes for machine learning*, (The MIT Press, 2005)
- T. Hastie, J. Friedman, and R. Tibshirani, *The Elements of Statistical Learning* (Springer, 2001)
- T. Mitchell, *Machine Learning* (McGraw-Hill, 1997)
- I. Witten and E. Frank, *Data mining – Practical Machine Learning Tools and Techniques* (2005)

Quiz



You want to machine-learn the HOMO energies of molecules.

Which family of machine learning methods would you use:

Unsupervised or Supervised?

If Unsupervised, would you use a Clustering or a Dimensionality Reduction algorithm?

If Supervised, would you use a Regression or a Classification algorithm?



You want to machine-learn the HOMO energies of molecules.

Which family of machine learning methods would you use:

Unsupervised or Supervised?

If Unsupervised, would you use a Clustering or a Dimensionality Reduction algorithm?

If Supervised, would you use a Regression or a Classification algorithm?

Supervised learning, Regression algorithm



You want to machine-learn the atomization energies of a set of molecules.
How do you need to describe your molecules for this purpose?



You want to machine-learn the atomization energies of a set of molecules.

How do you need to describe your molecules for this purpose?

Need to describe the molecules with molecular representations (and not with common molecular descriptors used in QSAR/QSPR approaches) that encode the basic information contained in the molecular Hamiltonian. For example: SLATM



A company provides you a list of 10000 molecules and asks to choose good catalysts for a given reaction.

You only have the budget for testing experimentally 1000 compounds.

Which steps would you follow?



A company provides you a list of 10000 molecules and asks to choose good catalysts for a given reaction.

You only have the budget for testing experimentally 1000 compounds.

Which steps would you follow?

First, we need to describe our molecules with molecular representations or descriptors.

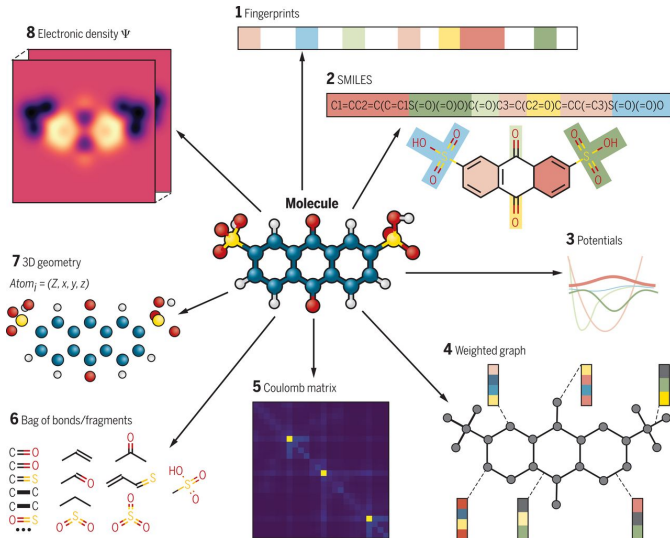
Then, we can use a clustering algorithm and pick the molecules from different clusters found to obtain a representative subset of 1000 molecules.

We run the experiments.

Then, we train a model on those 1000 compounds and predict the key property on the rest of the set.

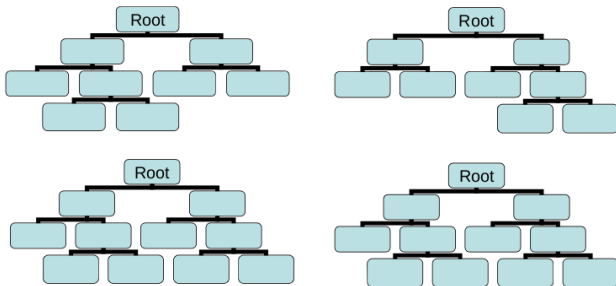
Appendix

Molecular representations: beyond this introduction





- Ensemble methods combine the predictions of several base estimators built with a given ML algorithm in order to improve generalizability / robustness over a single model
- Random forest algorithm is an ensemble technique that combines multiple decision trees
- Predictions are made by majority vote of the individual trees
- Better generalization than an individual decision tree due to randomness





- It runs efficiently on large databases
- It can handle thousands of input variables without variable deletion
- It gives estimates of what variables are important although less interpretative than single DT
- Performance increases with the number of trees until it saturates

