## Chimica Computazionale Machine learning and computational chemistry

### Emanuele Coccia



PhotoInduced Quantum Dynamics (PIQD) Group

< □ > < □ > < □ > < □ >

## Nobel prize in Physics 2024

## The Nobel Prize in Physics 2024

## John J. Hopfield

"for foundational discoveries and inventions that enable machine learning with artificial neural networks"

## **Geoffrey Hinton**

"for foundational discoveries and inventions that enable machine learning with artificial neural networks"





### The Nobel Prize in Chemistry 2024

David Baker

"for computational protein design"



### Demis Hassabis

"for protein structure prediction"



### John Jumper

"for protein structure prediction"

ヘロン 人間 とくほど 人間と



- Computer systems able of mimicking decision-making and problem-solving tasks of a human mind
- Machine learning:
  - Pathway to AI that uses statistical models and training algorithms
  - Learns insights and patterns in the data
  - Makes new predictions without additional input/programming
- Large amount of data available

イロン イ理シ イヨン イヨン

ML algorithms:

- Estimate relationships without any instruction of how to analyze or draw conclusions from the data
- Can recover mappings between a set of inputs/outputs or from the inputs alone
- Can discover structure in the data

< 口 > < 同 > < 臣 > < 臣 >

ML algorithms:

- Estimate relationships without any instruction of how to analyze or draw conclusions from the data
- Can recover mappings between a set of inputs/outputs or from the inputs alone
- Can discover structure in the data
- Use universal approximators

< 口 > < 同 > < 臣 > < 臣 >

- **Supervised Learning**: Learn from labeled data (regression, classification)
- **Unsupervised Learning**: Find patterns in unlabeled data (e.g., clustering, dimensionality reduction)
- **Reinforcement Learning**: Learn through reward-based interaction with the environment



イロト イポト イヨト イヨト



э

・ロン ・回 と ・ ヨ と ・ ヨ と

# Supervised Learning

- Input and output pairs (labeled data)
- Train model to learn mapping: f(x) = y
- Examples: from structure to spectrum, acid or base

# SUPERVISED LEARNING

Supervised machine learning is a branch of artificial intelligence that focuses on training models to make predictions or decisions based on labeled training data.



ヘロト ヘ回ト ヘヨト ヘヨト

# Unsupervised Learning

- No labeled output
- Aim: discover structure in data
- Examples: clustering from a MD trajectory



## **Reinforcement Learning**

- Agent interacts with environment
- Learns to maximize cumulative reward
- Examples: game playing, robotics



- Linear Regression
- Decision Trees
- k-Nearest Neighbors
- Support Vector Machines
- Neural Networks

< ロ > < 同 > < 三 >

< ≣ >

# What is a Neural Network?

- A neural network is a set of algorithms modeled after the human brain
- It is designed to recognize patterns from input data
- It consists of layers of interconnected neurons
  - Artificial neurons (nodes)
  - Weighted connections between neurons
  - Activation functions



- Input layer: Receives data
- **Hidden layers**: Process the data through weighted connections and activation functions
- Output layer: Produces final prediction or classification

## Artificial Neuron

Each neuron performs:

- Inputs: x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>
- Weights: *w*<sub>1</sub>, *w*<sub>2</sub>, ..., *w*<sub>n</sub>
- A weighted sum of its inputs:  $\sum w_i x_i + b$
- Output:  $y = f(\sum w_i x_i + b)$



2

イロン イ理 とくほ とくほ とう

## • Add non-linearity to the model



イロト イヨト イヨト イヨト

- Add non-linearity to the model
- Sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$
- ReLU: f(x) = max(0, x)
- Tanh:  $tanh(x) = \frac{e^x e^{-x}}{e^x + e^{-x}}$

- Add non-linearity to the model
- Sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$
- ReLU: f(x) = max(0, x)
- Tanh:  $tanh(x) = \frac{e^x e^{-x}}{e^x + e^{-x}}$

## Training a Neural Network

- Forward Propagation: Compute outputs layer by layer
- Loss Calculation: Measure prediction error
- Backward Propagation: Use gradients to adjust weights
- Optimization: Apply updates using algorithms like gradient descent
  - Common loss functions L:
    - Mean Squared Error
    - Cross-Entropy Loss
  - Weight update rule:

$$\mathsf{w} := \mathsf{w} - \eta \cdot \frac{\partial L}{\partial \mathsf{w}}$$

where  $\eta$  is the learning rate

< 口 > < 同 > < 臣 > < 臣 >

## • Gathering and preparing data

æ

イロト イヨト イヨト イヨト

- Gathering and preparing data
- Choosing a representation

イロト イヨト イヨト イヨト

- Gathering and preparing data
- Choosing a representation
- Training the model
  - Train model candidates
  - Evaluate model accuracy

< ロ > < 同 > < 三 >

- Gathering and preparing data
- Choosing a representation
- Training the model
  - Train model candidates
  - Evaluate model accuracy
- Testing the model out of sample

< D > < B > < E</p>

- Gathering and preparing data
- Choosing a representation
- Training the model
  - Train model candidates
  - Evaluate model accuracy
- Testing the model out of sample
- Deploy and monitor

< ロ > < 同 > < 三 >

- Quality and quantity of data
- Overfitting and underfitting
- Model interpretability
- Ethical and societal impacts

< ロ > < 同 > < 三 >

## Machine learning in a nutshell

#### **BIAS-VARIANCE TRADEOFF**

What is a good ML model?



#### **CROSS-VALIDATION**

How to find a good ML model?



## Machine learning in chemistry: overview



- Transformative impact on chemical sciences
- Dramatic acceleration of computations
- · Amplifying insights available from chemistry methods
- · Coaction of expertise in computer and physical/chemical sciences

< ロ > < 同 > < 三 >

## Machine learning in chemistry: overview

- ML = machine learning
- CC = computational chemistry
- CPI = chemical and physical intuition



Need robust data sets

Robust data sets

· Limited understanding

イロン イ理 とくほう くほう

## Catalyst accelerating data-driven hypotheses generation

## Machine learning in chemistry: overview

Occurrence of any ML term in American Chemical Society journals



22/34

## Data sets

- ML models as sophisticated parametrizations of data sets
- Data set must be representative
- Quality of data set determines the model effectiveness
- Avoid/reduce biases or artifacts
- CPI to reduce the function space
- A priori removing of unphysical solutions



Image: A matrix and a matrix

# CC databases for ML

| database                 | description   | location  |
|--------------------------|---|---|
| AFLOWLIB                 | databases containing calculated properties of over 625k materials <sup>510</sup>  | http://www.aflowlib.org   |
| ANI-1                    | large computational DFT database, which consists of more than 20 M off equilibrium conformations for 57.5k small organic molecules <sup>511,512</sup>   | https://github.com/isayev/ANI1_dataset  |
| ANI-1x/ANI-<br>1ccx      | ANI-1x contains multiple QM properties from 5 M DFT calculations, while ANI-1ccx contains 500k data points obtained with an accurate CCSD(T)/CBS extrapolation <sup>513</sup>                                   | https://github.com/aiqm/ANI1x_datasets  |
| BindingDB                | measured binding affinities focusing on interactions of proteins considered to be<br>candidates as drug-targets; 1 200 000 binding data for 5500 proteins and over 520 000<br>drug-like molecules <sup>11</sup> | http://www.bindingdb.org  |
| Clean Energy<br>Project  | contains ~10 000 000 molecular motifs of potential interest which cover small molecule organic photovoltaics and oligomer sequences for polymeric materials <sup>515</sup>                                      | http://cepdb.molecularspace.org   |
| CoRE MOF                 | database containing over 4700 porous structures of metal-organic frameworks with<br>publicly available atomic coordinates; includes important physical and chemical<br>properties <sup>16</sup>                 | 10.11578/1118280  |
| FreeSolv                 | experimental and calculated hydration free energies for neutral molecules in water <sup>517</sup>   | http://www.escholarship.org/uc/item/6sd403pz  |
| GDB                      | GDB-11, GDB-13, and GDB-17; together these databases contain billions of small organic molecules following simple chemical stability and synthetic feasibility rules <sup>518</sup>                             | http://gdb.unibe.ch/downloads/  |
| Hypothetical<br>Zeolites | contains approximately 1 M zeolite structures <sup>519</sup>  | http://www.hypotheticalzeolites.net/  |
| Materials<br>Project     | contains computed structural, electronic, and energetic data for over 500k $compounds^{520}$  | https://www.materialsproject.org  |
| MD17                     | data sets in this package range in size from 150k to nearly 1 M conformational geometries; all trajectories are calculated at a temperature of 500 K and a resolution of 0.5 fs <sup>372</sup>                  | http://www.sgdml.org  |
| MoleculeNet              | contains data on the properties of over 700k compounds <sup>521</sup>   | http://moleculenet.ai   |
| Open Catalyst<br>Project | 1.2 M molecular relaxations with results from over 250 M DFT calculations relevant for renewable energy storage $^{522}$  | https://opencatalystproject.org/index.html  |
| OQMD                     | consists of DFT predicted crystallographic parameters and formation energies for over 200k experimentally observed crystal structures <sup>523</sup>  | http://oqmd.org   |
| PubChemQC<br>PM6         | provides 221 million molecular structures optimized with the PM6 method and several electronic properties computed at the same level of theory <sup>524</sup>   | http://pubchemqc.riken.jp/pm6_datasets.html   |
| PubChemQC                | provides ${\sim}3$ million molecular structures optimized by DFT and excited states for over 2 million molecules using TD-DFT $^{525}$  | http://pubchemqc.riken.jp/  |
| QM7-X                    | comprehensive data set of 42 physicochemical properties for ~4.2 M equilibrium and nonequilibrium structures of small organic molecules with up to seven non-hydrogen (C, N, O, S, Cl) atoms <sup>226</sup>     | https://zenodo.org/record/4288677#.<br>X9jHNC2ZNTY  |
| QM9                      | geometric, energetic, electronic, and thermodynamic properties for 134k stable small<br>organic molecules out of GDB-17 <sup>527</sup>  | https://figshare.com/collections/Quantum_<br>chemistry_structures_and_properties_of_134_<br>kilo_molecules/978904 |
| Synthesis<br>Project     | collection of aggregated synthesis parameters computed using the text contained within over 640 000 journal articles $^{528}$   | www.synthesisproject.org  |
| quantum-<br>machine.org  | a repository of diverse data sets, including valence electron densities, chemical reactions,<br>solvated protein fragments, and molecular Hamiltonians  | http://quantum-machine.org/datasets/  |

## Benchmarking data sets

- Learning curves (QM9 database, 134K organic molecules)
- Target and DFT/B3LYP accuracy



E. Coccia (DSCF)

- Searching stationary points of a PES
- Generating force fields for MM and MD
- Use in metadynamics (collective variables)
- Chemometrics
- Text mining for extracting scientific information
- Structure/property relationship in spectroscopies
- Retrosynthesis
- Materials
- Drug design
- ...

イロン イ理シ イヨン イヨン

- One-to-one spectrum-structure relationships
- Conventionally with CC

- One-to-one spectrum-structure relationships
- Conventionally with CC
- Machine-learning protocol to correlate spectral fingerprints with local molecular structures
  - Quick and accurate prediction of infrared (IR) and Raman spectra
  - Structure recognition of functional groups from vibrational spectral features

イロト イヨト イヨト -

# IR and Raman with quantum chemistry

## Vibrational modes

- Diagonalization of the mass-weighted Hessian matrix
- Eigenvectors: normal modes **q**
- Eigenvalues: frequencies
- Harmonic approximation

# IR and Raman with quantum chemistry

## Vibrational modes

- Diagonalization of the mass-weighted Hessian matrix
- Eigenvectors: normal modes **q**
- Eigenvalues: frequencies
- Harmonic approximation

• IR

• Change in the dipole moment  $\mu$ 

IR intensity 
$$\propto \left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\mathsf{q}}}\right)^2$$

< 口 > < 同 > < 臣 > < 臣 >

# IR and Raman with quantum chemistry

## Vibrational modes

- Diagonalization of the mass-weighted Hessian matrix
- Eigenvectors: normal modes **q**
- Eigenvalues: frequencies
- Harmonic approximation

IR

• Change in the dipole moment  $\mu$ 

IR intensity 
$$\propto \left(rac{\partial oldsymbol{\mu}}{\partial oldsymbol{q}}
ight)^2$$

- Raman
  - Change in the polarizability lpha

Raman intensity 
$$\propto \left(rac{\partial oldsymbol{lpha}}{\partial oldsymbol{q}}
ight)^2$$

イロト イポト イヨト イヨト

## Applications: vibrational spectroscopy

- Hydroxyl (OH, 3000-4000 cm<sup>-1</sup>) and carbonyl (C=O, 1400-2000 cm<sup>-1</sup>) groups
- Dataset with around 21,000 molecules
- Spectra with DFT/ B3LYP/6-31G(2df,p)



 $MAE = \frac{1}{n} \sum_{i} |x_i - x_{\text{ref},i}|$  $MRE = \frac{1}{n} \sum_{i} |x_i - x_{\text{ref},i}| / |x_{\text{ref}}|$ 

E. Coccia (DSCF)

ヘロン 人間 とくほ とくほ とう

## Applications: vibrational spectroscopy

## Histidine



E. Coccia (DSCF)

# Applications: molecular and material design

- Identify compounds with desired properties (high-throughput screening)
- OLED emitters (*k*<sub>TADF</sub> delayed fluorescence rate constant)
- ML comparable to CC calculations, at a fraction of the computational cost



# Applications: retrosynthesis

- Design of chemical steps
- SCScore: data-driven metric specific for reactions
- Monotonic increase in complexity with SCScore
- ML to overcome the generalization issues of rule-based algorithms
- Synthesis of a precursor to lenvatinib



# Applications: catalysis

Accelerated discovery of CO2 electrocatalysts using ML



E. Coccia (DSCF)

# Applications: drug design

### Discoidin domain receptor 1 (DDR1)



E. Coccia (DSCF)