

Machine Learning techniques for predicting molecular properties

Viviana Petrescu
I&C, EPFL

Abstract—The high computational cost of quantum chemistry calculations have prompted the use of less expensive machine learning methods for predicting molecular properties in chemical compound space. Finding good feature representations for molecules is hard, in part because of the graph-like structure geometry of the molecules that need to be represented as high dimensional vectors. [Add what this paper is about](#)

Index Terms—thesis proposal, candidacy exam write-up, EDIC, EPFL

I. INTRODUCTION

The discovery of new molecular materials in chemistry has the potential of solving many of the problems we face today. Having a system which predicts both accurately and at a small computational cost the properties of new materials is highly desirable and has applications ranging from novel drugs discovery, water purification to efficient materials design for high energy transmission and storage [1].

II. BACKGROUND WORK

[This write-up serves two purposes. First, it forms the basis for your candidacy exam. As such you should summa-](#)

Proposal submitted to committee: June 13th, 2009; Candidacy exam date: June 20th, 2009; Candidacy exam committee: Exam president, thesis director, co-examiner.

This research plan has been approved:

Date: _____

Doctoral candidate: _____
(name and signature)

Thesis director: _____
(name and signature)

Thesis co-director: _____
(if applicable) (name and signature)

Doct. prog. director: _____
(B. Falsafi) (signature)

size the three papers selected by your advisor and yourself, and analyze as well as discuss them critically. Second, the write-up is also your thesis proposal. Therefore, the last one or two pages should be dedicated to your own preliminary work. A road-map of how you plan to advance the state of the art in your chosen area should also be given. For further details please consult the document “PhD Candidacy Exam Overview.” You can find the latest version at <http://phd.epfl.ch/page57746-en.html>. Describe briefly the context, the problem, shortcomings in prior approaches, and your proposed approach and solution. Forecast results. Background — Describe the three papers in detail, the problem they tackle, the solutions and results, and their shortcomings, and how they relate to your work. This part builds the basis for the oral candidacy exam.

A. Learning Invariant Representations of Molecules for Atomization Energy prediction

Representing Molecules

While domain specific descriptors for molecules exist [5], recent work [4] has proposed to predict properties of a molecule of size N only from the 3D positions of the atoms R_i $i \in 1..N$ and their nuclear charge Z_i $i \in 1..N$. This has prompted the introduction of the Coloumb Matrix descriptor, whose individual entries appear in the Schroedinger equation. It is defined by a $N \times N$ matrix with entries given by:

$$M(i, j) = \begin{cases} 0.5 * Z_i^{2.4} & \text{if } i = j \\ Z_i * Z_j / |R_i - R_j| & \text{otherwise} \end{cases} \quad (1)$$

where $R_i - R_j$ represents the distance between the atoms i and j .

The dimensionality of the Coloumb Matrix is given by the number of atoms in a molecule. Since that varies across molecules in a dataset, one common trick is to pad with 0's the matrices corresponding to small molecules until they reach the maximum molecule size in a dataset. This limits the size of the molecules that can be used, since the descriptor has complexity $O(N^2)$, where N is the number of atoms.

Desired properties of descriptor Due to the graph-like structure of the molecules, finding a fixed size representation is difficult. The desired properties of a molecular descriptor are invariance to translation and rotation of the molecule and invariance to the indexing of the atoms.

Solved using sorted or Random Coloumb While the Coloumb Matrix representations solves the rotation and translation invariance through the use of distances between atoms

$R_i - R_j$, invariance in atom indexing still needs to be tackled since any permutation of atom indexes results in a valid Coloumb descriptor. Two variations of the descriptor are proposed in [3]. The first one, called Sorted Coloumb, uses the permutation of the atoms given by the sorting of the row norms of a valid Coloumb matrix. Any molecule has a unique representation given by the Sorted Coloumb matrix. The second representations is based on the idea that the norms of the rows can have very similar values in practice and the sorting, therefore also the atom indexing, is subject to small noise. The new descriptor, Random Coloumb Matrices is a collection of Sorted Coloumb matrices. For every molecule, approximately 10 Randomly Sorted Coloumb matrices are drawn. One Random Sorted Coloumb Marix is obtained by sorting the set of rows was sorted according to their norm at which a small Gaussian noise was added. New predictions can be made by averging the prediction results.

Reach state of the art result at the time. The current new value is 1.5 using Bag of Bonds. The performane of the descriptors was evaluated on predicting the atomization energy on a dataset of 7165 molecules with at most 23 atoms per atom. The predicted values range from -800 to -2000 kcal/mol. The splitting of the data into 5 folds for cross validation is done using stratified sampling. The molecules were clustered into 5 sets with similar atomization energy levels and the folds were created by randomly selecting one molecule from each bucket.

While this is a good strategy for obtaining good generalization performance, in normal machine learning practice we should not touch the test dataset. If we test the performance on a new molecule whose atomization energy is not within the bounds present in the training dataset, it is likely that the prediction will not be very accurate. Moreover, if we have a molecule with more atoms than present in the dataset, we would need to retrain our model with a descriptor of diferent size. The same problem appears across chemical compound space if we try to predict the atomization energy of a molecule with the same number of atoms but with different atom types composition. In general, we do not have guarantees that the a molecule with similar atom composition and similar atom type have target prediction in the same range as our training set.

For experimental evaluation of the performance of descriptors, both kernel ridge regression and multy layer feed forwards networks were tried. The best test performance of 3.1 kcal/mol MAE was obtained using a neural network with Randomly Sorted Coloumb matrix. The improvement was three fold with respect to the previous state of the art, which gave a MAE of 9kcal/mol. A performance level is 1kcal/mol to reach chemical accuracy. [Maybe put this as footnote](#) Later work [2], improved upon this result by reducing the current state-of-the-art to only 1.5kcal/mol using a bag of bonds descriptor with Laplacian kernel.

Difficulty of training a neural network, maybe cite tricks of the trade While recently Neural Nets (especially deep networks) have achieved state-of-the-art in many fields ranging from computer vision, natural language processing and speech recognition, one of the main challenges they pose is the

difficulty in training for people without a lot of experiences with nnets in particular.

Cite here Neural nets tricks fo the trade In order to obtain the competitive performance of 3.1kcal/mol using neural nets, multiple tricks have been used, inspired from cite here NN. First of all, taking the real entries of the Coloumb matrix as input to the neural net proved to perform poorly. For this, a binarization step was performed which added an extra dimensionality to the dataset. Specifically, if $x \in R^D$ represents a flatten Coloumb matrix, it is mapped to a new binarized descriptor $y \in R^{D \times M}$ such that $y_i \in R^{1 \times M}$ is computed by taking shifted versions of x_i entry that are passed through a sigmoidal function

$$y_i = [\dots, \tanh(\frac{x_i - \theta}{\theta}), \tanh(\frac{x_i}{\theta}), \tanh(\frac{x_i + \theta}{\theta})] \quad (2)$$

Depending on the range of every x_i , the size of y_i can vary acrosss dimensions if we ignore saturated values. Other parameters that are selected using cross validation are the number of hidden units per layer, the number of layers, the learning rate.

B. Self-Taught Learning: Transfer Learning from Unlabeled Data

Present the different machine learning framewroks that exists, with picture. It is a two part optimization problem. - First find bases - Then find new activations for the old ones. Through experimental results - since they were tested on a range of problems from image, text and sound.

C. Information-Theoretic Regret Bounds for Gaussian Process Optimization in the Bandit Setting

Optimizing an unknown function is an important problem.

Optimization of a black box function that is expensive to evaluate is a common problem with applications ranging from active user modelling, hierarchical reinforcement learning and more recently, hyper parameter optimization of machine learning models. This problem can be posed as a multi-armed bandit problem where the function to be estimated is either drawn from a Gaussian Process or has low norm in reproducing Kernel Hilbert space.

Often done using Gaussian processes in a Bayesian setting. Many heuristic have been proposed for defining an acquisition function. Few is known for its conergence.

Among the common used acquisition functions we note the probabily of improvement, maximum expected improvement or upper confidence bound acquisition, with the latter two experimentally shown to perform better in practice. Besides their practical success, few is known in practice about their convergence properties for finite and infinte input dimensional space.

Sublinear converge rates are proven for GP-UCP. The bound is in two steps. First bound it on the information gain, then on the spectral of the kernel matrix which is further bound on the spectral operator.

Since RKHS is not relevant to our current setup, it will be skiped from the previous explanation.

1) *Gaussian Processes basics:* A Gaussian Process (GP) defines a distribution over functions f , $f \sim GP(u(x), k(x, x'))$, in a similar way in which probability densities define distributions over random variables. It is fully specified by its mean and its covariance matrix. In a Bayesian setting, we assume our function f is sampled from a Gaussian with prior $GP(0, k(x, x'))$, where k is the kernel or covariance matrix. The most common kernel types are the linear, squared exponential and Matern kernel. For T sampled points with noise, we have $y_t = f(x_t) + \epsilon_t$, where $\epsilon_t \in N(0, \sigma^2)$. The posterior over f , conditioned on the $y_{1:T}$ and $x_{1:T}$ in the presence of noise is given by $P(f_T | y_{1:T}, x_{1:T}) = N(u_T(x), \sigma_T^2(x))$, where:

- K_T is positive definite with entries $k(x, x'), x, x' \in A_T$
- $k_T(x) = [k(x_1, x), \dots, k(x_T, x)]^T$
- $k_T(x, x') = k(x, x') - k_T(x)^T (K_T + \sigma^2 I)^{-1} k_T(x')$
- $u_T(x) = k_T(x)^T (K_T + \sigma^2 I)^{-1} y_T$
- $\sigma^2(x) = k_T(x, x)$

2) *Information Gain and Experimental Design:* Unlike function optimization where the goal is to find the maximum of a function, in experimental design the goal is to find a good approximation of the function globally with few samples as possible. We note that the same algorithm can not be employed for both tasks, since in function optimization we might want to avoid sampling in certain regions where we are confident that the function has very small values.

For a sampled set A , we define $y_A = f_A + \epsilon_A$, $f_A = [f(x)]_{x \in A}$ and the noise $\epsilon_A \sim N(0, \text{diag}(\sigma^2))$. The information gain as

$$I(y_A; f) = H(y_A) - H(y_A | f) \quad (3)$$

where $H(y_A)$ and $H(y_A | f)$ represent the entropy or uncertainty in y_A and the uncertainty in y_A if we know f . In other words, information gain expresses the reduction in uncertainty about f given y_A . **Shouldn't be the other way around?**

Using the fact that for a Gaussian $H(N(u, \sigma)) = \frac{1}{2} \log |I + \sigma^{-2} K_A|$ we can rewrite $I(y_A; f) = I(y_A; f_A) = \frac{1}{2} \log |2\pi e \Sigma|$. Finding the subset A in the domain D with $|A| \leq |T|$ that maximizes IG is NP-hard. In practice, greedy approximation solutions are used that find a near-optimal solution by choosing sequentially points in A with maximum variance:

$$x_t = \arg \max_{x \in D} \sigma_{t-1}(x) \iff x_t = \arg \max_{x \in D} F(A_{t-1} \cup x) \quad (4)$$

with $t \in 1 \dots T$

The results from cite1 show that $F(A)$ is a submodular function and thus it holds that

$$F(A_T) \geq (1 - \frac{1}{e}) \max_{|A| \leq T} (F(A)) \quad (5)$$

From the above equation we see that our greedy choice of A_T leads to a near optimal solution.

3) *GP-UCB: maybe introduce some* For function optimization a good sampling sequence x_t is one that makes a tradeoff between exploration and exploitation. By choosing points with high variance (exploration) we try to reduce the overall uncertainty. By choosing points with high mean (exploitation) we try to concentrate samples in a region where we are more confident that is close to the optimal. This idea is expressed by choosing

Kernel	Linear	RBF	Matern
Regret R_T	$d\sqrt{T}$	$\sqrt{T(\log(T)^{d+1})}$	$v + d(d+1)/2v + d(d+1)$

TABLE I
REGRET BOUNDS.

in every round t , a point $x_t = \arg \max u_{t-1} + \beta_t^{1/2} \sigma_{t-1}(x)$ with β_t defined in [1]. The function that is being maximized is also called the acquisition function.

To sum up, at every round, the GP-UCB algorithm selects a point according to the above equation, samples $y_t = f(x_t) + \epsilon$ and updates u_t and σ_t using Bayes rule.

An important property of the acquisition function is that as defined by equation above is that it does not choose points x for which the upper confidence value is smaller than the maximum lower confidence value found so far. This can prune a large subareas of the search space.

4) *Regret Bounds:* The effectiveness of a sampling scheme in finding x^* (may not be unique) s.t $x^* = \arg \max_{x \in D} f(x)$ is evaluated using the *instantaneous* and the *cumulative regret*. The instantaneous regret $r_t = f(x^*) - f(x_t)$ is the error that we make at round t for choosing point x_t instead of x^* . The cumulative regret is the regret accumulated up to this round $R_T = \sum r_t$.

This paper cite here provides convergence rates for the the cumulative regret R_T in both the finite and infinite case. They show that $\lim_{T \rightarrow \infty} \frac{R_T}{T} = 0$ and that . For the most popular kernels, this is shown to be

5) *Bounding the regret using Information Gain:* The first step of the proof is to bound R_T by the maximum information gain $\gamma_T = \max_{A \subset D, |A|=T} I(y_A; f_A)$. In particular, in citep it is shown that with high probability

- $Pr(R_T \leq \sqrt{C_1 T \beta \gamma}) \leq 1 - \delta$ for finite D
- $Pr(R_T \leq \sqrt{C_2 T \beta \gamma} + 2) \leq 1 - \delta$ for D compact and convex.

where $\delta \in (0, 1)$, appropriate constants C_1, C_2, β_t variable from the acquisition function. The latter proof makes further assumptions that the function is smooth, which says that the probability of the partial derivatives of the function to have a large value is small.

6) *Bounding the Information Gain using the empirical spectrum:* The case of infinite input dimensionality are treated similarly with the finite case by assuming the existence of a discretized set $D_T \subset D, T \in \mathbb{N}$ with nearest neighbor distance $O(1)$, dense in the limit. The value of $\gamma_T^* = \max_{A \subset D_T, |A|=T} I(y_A; f_A)$ can be bounded by the greedy maximization value since it is a submodular function. By using another relaxation procedure, this value is further bounded by an expression depending on the eigenvalues λ_t of the kernel matrix $K_{D_T} = k(x, x')_{x, x' \in D_T}$ as follows

$$\gamma_T^* \leq \frac{0.5}{1 - e^{-1}} \max_{m_1, m_2, \dots, m_T} \sum \log(1 + \sigma^{-2} m_t \lambda_t^*) \quad (6)$$

where $m_t \in \mathbb{N}$, $\sum_t m_t = T$ and $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots$. Here we have considered a worst case scenario of assigning eigenvectors of K_T

In the next step, the information gain is bounded using the tail spectrum of K_{D_T} , which is $B(T^*) = \sum \lambda_t^*$ and $n_T =$

$\sum \lambda_t^\sim$ for any $T_* = 1..T$

$$\lambda_T \leq O(\sigma[B(T_*T + T_* \log n_T T)]) \quad (7)$$

From equation above, we conclude that if the tail spectrum $B(T_*)$ is small, or in other words, the eigenspectrum of K_{D_T} decays fast, the bound on γ_T is more tight.

7) *Bound the empirical spectrum by the operator spectrum:* For a kernel k with $k(x, x)u(x) = 1$ where $u(x)$ is uniformly distributed on D , Mercer's theorem the existence of a discrete eigenspectrum $(\lambda_s(x), \phi_s(x))$ with $k(x, x') = \sum \lambda_s \phi_s(x) \phi_s(x')$

The contribution bounds the empirical tails spectrum $B(T_*)$ by the kernel operator tails eigenspectrum $B_k(T_*)$. Analytical expression bounds on $B_k(T_*)$ for the most common kernels are given by Seeger.

8) *Conclusions:* The major contribution of the paper is making a connection between experimental design and GP optimization which leads to sublinear regret bounds for GP-UCB. The bound is derived in two steps. First the cumulative regret is bounded by the maximum information gain. In a second step, the max information gain is bounded by the tails of the eigenspectrum which in turn is bound by the tails of the kernel operator spectrum. The bound is tighter, the more rapidly the kernel spectrum decays. Its major importance is also

Although convergence rates are proven for infinite input dimensional spaces, the complexity of the GP step that incorporates the inversion of the covariance matrix makes it cubic in the number of samples T , which makes it often impractical for high dimensional input.

III. RESEARCH PROPOSAL

A.

Proposing a new descriptor which is invariant to permutations of the atom. 2.6+/- or 2.1 with noise. (still gives comparable accuracy to same DFT models)

B.

Propose of augmenting the data sets more easily and make the cross validation less sensitive to splitting -at the moment: take non H atoms, then sort then do CV.

Pose the problem as one of the semi supervised, self taught learning etc problems and try to make it generalize across compound space, or learn new embeddings of the atoms.

Although the use of Gaussian processes in material design is not new, its major drawback is the computational bottleneck. In our scenario, this can be used for tuning the hyperparameters of the model trained. Here the input dimensionality is given by the number of the hyper parameters (learning rate, activation function, number of hidden layers) used and the function to be minimized is the cross validation error. The result presented in the previous subsection were obtained like that. Talk if we have time about bayesian neural nets, were simpler models

outperform more easy models just by using hyper parameter optimization instead of grid search.

Write how you propose to advance the state of the art given the background. What is new technically? How does it improve over prior work? Summarize, suggest an approximate timeline, and list references.

REFERENCES

- [1] Johannes Hachmann, Roberto Olivares-Amaya, Sule Atahan-Evrenk, Carlos Amador-Bedolla, Roel S. Sanchez-Carrera, Aryeh Gold-Parker, Leslie Vogt, Anna M. Brockway, and Aln Aspuru-Guzik. The harvard clean energy project: Large-scale computational screening and design of organic photovoltaics on the world community grid. *The Journal of Physical Chemistry Letters*, 2(17):2241–2251, 2011.
- [2] O.A. von Lilienfeld K.-R. Müller K. Hansen, F. Biegler and A. Tkatchenko. *Interaction Potentials in Molecules and Non-Local Information in Chemical Space*. Phys. Rev. Lett. (March 10, 2014)., 2014.
- [3] Grégoire Montavon, Katja Hansen, Siamac Fazli, Matthias Rupp, Franziska Biegler, Andreas Ziehe, Alexandre Tkatchenko, Anatole V Lilienfeld, and Klaus-Robert Müller. Learning invariant representations of molecules for atomization energy prediction. In *Advances in Neural Information Processing Systems*, pages 440–448, 2012.
- [4] Matthias Rupp, Alexandre Tkatchenko, Klaus-Robert Müller, and O. Anatole von Lilienfeld. Fast and accurate modeling of molecular atomization energies with machine learning. *Phys. Rev. Lett.*, 108:058301, Jan 2012.
- [5] Roberto Todeschini and Viviana Consonni. *Handbook of molecular descriptors*, volume 11. Wiley-VCH, 2000.