

UiO : **Department of Physics**
University of Oslo

Sondres master utkast

16.01.2020

Sondre Torp - Sondrt@student.matnat.uio.no



Abstract

This is the Abstract!

Acknowledgements

I like to acknowledge ... Sabrina!

UOC

Georges

m.m.

Contents

I	Introduction	1
1	Introduction and Overview	1
1.1	Motivation	1
1.2	Scope of the thesis	1
1.2.1	Research Question	1
1.2.2	Approach	1
1.3	Structure of the thesis	1
II	Foundations	2
2	Batteries	2
2.1	Cell operation principles	2
2.2	Cell limitations	4
2.2.1	Polarization	4
2.2.2	Rate-capacity effect	4
2.2.3	degradation/aging	4
2.3	Battery chemistries	4
2.4	Intercalation batteries	4
2.5	Electrodes and features	4
3	Machine Learning	8
3.1	The basics of Machine Learning	8
3.1.1	Example time!	8
3.1.2	Supervised and Unsupervised Learning	9

3.1.3	Regression and Classification Problems	9
3.1.4	Features and Feature Selection	10
3.2	Random Forest	10
3.3	Support Vector Regression	10
3.3.1	Radial Basis Function Kernel	10
3.3.2	SVR and the Bias-Variance-Trade-off	10
3.4	Mean square error and Root mean square deviation	10
3.5	R^2 score - The Coefficient of Variation	11
3.6	Principle Component Analysis	11
3.7	K-fold cross validation	12
3.8	Mutual Information**	12
3.9	Machine Learning \times batteries	12

III Method 13

4 Some Method title 13

4.1	Data set and Experimental Environment	13
4.2	Volumetric number density	14
4.3	Void Fraction	14
4.4	AP-RDF Descriptors of Electrode materials	15
4.5	Algorithm	16

IV Result & Discussion 19

5 Result section title 19

5.1	Random factors from database.	19
5.1.1	Average Voltage	19

5.1.2	Capacity	19
5.1.3	Energy Density	19
5.2	Volumetric number density	19
5.2.1	Average Voltage	19
5.2.2	Capacity	20
5.2.3	Energy Density	20
5.3	Void fraction	20
5.4	AP-RDF	20
5.5	Stability	20
5.6	Geometrical descriptors	20

V Summary 21

6 Summary and future work 21

6.1	Batteries	21
6.2	future work	21
6.2.1	improving method	21

List of Figures

1	Wall of numbers.	19
---	--------------------------	----

List of Tables

Part I

Introduction

1 Introduction and Overview

1.1 Motivation

I want a job, or something.

1.2 Scope of the thesis

Clarify what you want to cover.

Batteries are vastly complex and much efforts have been devoted to the development of these. Yet, with all these efforts, it still is a never ending chase for batteries that can push the limits of their properties even further. This work proposes a methodology to predict these properties accurately without the need of big scale simulations, or computer heavy calculations. Using state of the art machine learning, and base properties of all ready existing databases, we propose a set of predictors to see if we can predict the properties of new, undiscovered electrodes, or even new properties in already well known electrodes.

1.2.1 Research Question

How to better batteries?

RQ1: Is there potential for the use machine learning to ease the search for better battery materials?

RQ2: Which ML method would be the most optimal for such a search?

RQ3: What predictors are the most suited for such a task, and which would yield the most efficient training.

RQ4: Does the size of the database matter? Or is it possible to find a solution with a "not optimal", or even small, database?

1.2.2 Approach

The choice of features examined in this work is inspired by an extensive survey done on similar project especially in the field of metal organic framework done by my supervisors in Crete, and dictated, to some degree by the lack of more data.

In order to evaluate the effect of different features, a prediction approach using principle component analysis was utilized. First we decided to use physical descriptors, that is for instance geometrical properties of the unit cell. This because it was greatly efficient in other studies(REF), and is straightforward.

Then we needed to find other descriptors, and the void fraction seemed like the next obvious thing

1.3 Structure of the thesis

First the most essential concepts from the fields of batteries, machine learning, and work already explored on these two fields conjoined, are introduced. Then the method will be explained before rounding up our results so far before trying to put this all into perspective.

Part II

Foundations

2 Batteries

Before doing anything else, the basis of batteries needs an introduction. In this section, this will be done, and some of the most essential properties, in the writers opinion, will be introduced.

2.1 Cell operation principles

Electrochemical cell and basic concepts

A magnesium- or lithium-ion cell or battery consists of a positive and a negative electrode(s) and a electrolyte in a casing. The electrodes function as active materials which can accept or release Mg or Li ions; a conductive additive which electrically connects the active material with a current collector; and a suitable binder which attaches the electrode particles to the current collectors. The current collector enable connection to an external circuit. The cell will normally also include a separator which usually is a semi-permeable membrane that is situated between the electrodes. The separator permits ionic charge carriers to travel through the electrolyte from one electrode to the other while separating the electrodes.

While charging or discharging the battery over the external circuit, Mg- and Li- ions move between the positive and negative electrodes (IMAGE TIME!)- When the cell is discharge, electrochemical reduction takes place at the positive electrode as electrons flow through an external electrical load towards the positive electrode while cations move within the cell through the electrolyte to the positive electrode. At the negative electrode, oxidation occurs. The positive and negative electrode is commonly referred to as the *cathode* and *anode*, respectively. These two materials have different voltage, high and low. This difference is the cell voltage which is the driving force for the discharge of the battery. For secondary batteries, as are the one in question, it is possible to recharge batteries by reversing this process by applying an external electrical power source which applies an over-potential, that is a higher voltage than the one produced by the cell, with the same polarity. This reverses the movements of cations and electrons.

Repeating some of this in the capacity section. Check this later.

The capacity can be defined as:

$$C = \int I(t) \cdot dt \quad (2.1)$$

And is the i number of electrons or cations exchanged between the negative and positive electrodes. $I(t)$ is the current, i.e. the number of electrons flowing over the external circuit per time

interval dt which is integrated over the discharge period. The capacity is normally expressed as Ah/kg. The battery can deliver a power that is defined as

$$P(t) = V(t)I(t) \quad (2.2)$$

Where $I(t)$ is the current, defined as earlier, drawn at a cell voltage $V(t)$. The amount of work that can be done by the battery, or the energy contained in the battery, is then defined as the power delivered over the discharge period

$$W = \int P(t) \cdot dt = \int V(t)I(t) \cdot dt \quad (2.3)$$

Specific capacity and energy densities of battery materials can be compared relative to mass, volume and cost. The more electrode material that a battery contains, the greater is its capacity and energy. The higher the cell voltage the greater its power and energy.

The active materials of the electrodes allow the reversible uptake and release of Mg, or Li ions. This may happen by; movement of the Li, or Mg ions into, i.e. *intercalation* or *intercalation* or out of, i.e. *extraction* or *deintercalation*, their chemical structures, *phases*, by conversion of the materials between Li/Mg poor and rich i.e. *alloying* or rich and poor, e.g. *dealloying* phases, or by conversion of the electrode material into other more Li/Mg rich/poor chemical forms or mixtures, usually referred to as *conversion* or *displacement* reaction, with the average Li/Mg content of the entire electrode varying. The total Li or Mg content in the electrodes will thus either be varied by changing the composition of one phase or the ratio between coexisting phases. In this work ?we will only look at *intercalation* type batteries, due to the *database*, more on this later.

2.2 Cell limitations

2.2.1 Polarization

2.2.2 Rate-capacity effect

2.2.3 degradation/aging

2.3 Battery chemistries

2.4 Intercalation batteries

2.5 Electrodes and features

In this section the features used in ML as predictors will be introduced. First will the pair properties be introduced, before going into the more electrode specific features.

As a general note. These features are based on optimal design and discharge conditions. These values are helpful to set a number on the "goodness" of a battery, the actual performance may vary under normal conditions of use. **Nice to give this note?**

Average Voltage

The theoretical voltage and capacity of a cell are function of the anode and cathode materials, with the composition of the electrolyte, and the temperature, normally 25°C.

The active materials contained in the cell determines the standard potential, E^0 , which can be calculated from the free-energy. The standard potential of a cell can be calculated from the standard electrode potential:

$$\text{Anode(oxidation potential)} + \text{cathode (reduction potential)} = \text{standard cell potential} \quad (2.4)$$

The cell voltage is also dependent on other factors including concentration and temperature, as expressed on the nernst equation. (REF) *Average Voltage* as we use, is defined as the voltage average during the discharge. It is lower than the theoretical voltage.

Capacity

Capacity represents specific energy in Ampere-hours(Ah), and is the discharge current a battery can deliver over time. The capacity is also determined by the amount of active materials in the cell, expressed through the total quantity of electricity involved in the electrochemical reaction and is defined in terms of coulombs or ampere-hours. Theoretically, 1 gram equivalent weights of the active material in grams divided

Specific Energy

Specific Energy, or gravimetric energy density, defines battery capacity in weight, energy density, or volumetric energy density, defined as:

$$\text{Watt-hours/gram} = \text{Voltage} \times \text{Ampere-hours/gram} \quad (2.5)$$

Physical stability

What we refer to as Physical stability is Energy above hull. The energy that is demanded for decomposition of the material into the set of most stable materials at that chemical composition. Some positive value indicates that the material is not stable. While a zero energy above hull indicates that this is the most stable material at its composition.

Cycle life

Rate capability

RC

Self discharge

SD

Energy per atom

EpA

volume

Volume of the unit cell defined as **This is too dumb? No, Explain what type of volume**

Formation energy per atom

F_{epa}

Band gap

The band gaps of a solid is simply the range of energies an electrode in a solid can not have. While the bandstructures. **How much to include? Should I here have a page on quantum physics and the bandstructure? - Only include relevant stuff.**

Total magnetization

T_m

Elasticity

E

Porous Electrodes

In a fuel cell system, the reactant is supplied from the electrolyte phase to the catalytic electrode surface. Electrodes are often composites made of active reactants, binders and fillers, in batteries. To minimize the energy loss of both activation and concentration polarizations at the electrode surface and to increase the electrode efficiency or utilization, it is often preferred to have a large electrode surface area. This can be done by having a porous electrode design. A porous design can provide an interfacial area per unit volume that is considerably higher than that of a planar electrode. A porous electrode is an electrode that consists of a porous matrix of solids and void space. The electrolyte penetrates the void space of a porous matrix. In such an active porous mass, the mass transfer condition in conjunction with the electrochemical reaction occurring at the interface is very complicated. In a given time during cell operation, the rate of reaction within the pores may vary significantly depending on the location. The distribution of current density within the porous electrode depends on the physical structure (pore size), the

conductivity of the solid matrix and the electrolyte, and the electrochemical kinetic parameters of the electrochemical processes.

3 Machine Learning

In this chapter we summarize some concepts of machine learning and related ideas. The first section introduces the basic ideas behind machine learning and *some of the best known examples* will be presented. Secondly the concepts of supervised and unsupervised learning will be presented with a clarification on the difference between regression and classification problems, so that we can define where in the field of machine learning this work resides in. Before we round off this section with a brief explanation on the role of data, how features can affect the effectiveness of a model, and finalizing with the concepts of over- and under-fitting, and how these are related to the bias-variance-trade-off.

The following section explains the basics of methods utilized in this work, starting by first giving an introduction to Random forest and support vector regression(SVR). Subsequently a short description of the validation methods used in this work is given. These are; K-fold cross validation and how it is used in optimizing our random forest method and mean square error(MSE).

Sondre: Did you forget something? Come back to this when done with the section.

3.1 The basics of Machine Learning

Machine learning comes from the field of pattern recognition and learning theory, and is defined as the field of study that gives computers the ability to learn without being explicitly programmed. Or more precise: "... A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with the experience E ..."(**mitchell1997machine**). At its core the ability to learn by detecting patterns in usually huge amounts of data that, more often then not, is impossible to perceive for a human.

3.1.1 Example time!

As an introduction on how machine learning was applied to learn and recognize patterns in our work, it will be useful to start with a simple example applied to the recognition of the handwritten number "5". (PICTURE TIME!)

How two people writes a single digit may vary to an extensive degree. It might seem to be a easy problem, but if the recognition is to be done manually million of times, it is no longer a trivial task for any one human being. Therefore a model which can recognize these digits would be useful. A model that takes a picture of a digit and outputs that digit in a way that is recognizable for a machine, that is, a digital format.

Machine learning only works when you have data, preferably a large amount of data. For instance data from the MNIST test dataset(REF). This database contain 60,000 images of hand-

written numbers for this very purpose. The images all are 18x18 pixels. The data is divided into two sets, one training set: X_{Train} and one test set: X_{test}

How do one represent an image as something that makes logical sense to a computer? Most learning algorithms take numbers as input. To a computer one image is nothing more then a grid of numbers that represent how dark a pixel is. So each picture contains a gray-scale value that ranges from 0 to 255. Where each sample can be viewed as a vector consisting of 324 *features*. Every sample has a corresponding label value, or *target*, which is the digital equivalent to the handwritten sample. We let the corresponding targets be denoted: Y_{train} and Y_{test} . Next we designate our *learner* denoted by function A . A is then given our training set S , where $S = (X_{train1}, Y_{train1}), \dots, (X_{trainN}, Y_{trainN})$ and returns a prediction rule: $h : X \rightarrow Y$. This rule is also called a predictor, in general, a classifier, or a regressor, depending on the problem in question.

The *training phase* is a process where the learning algorithm gets tweaked to best capture the correlating structure of the data set, so that it can better predict new data. As mentioned in the last paragraph the output from the *training phase* is called a *predictor*. The next step is to introduce the *predictor* for new, unseen data, so that it can be classified. Then we compare the Y_{test} to our predicted value Y_{pred} given by h to see if our model generalizes well to unseen data in X_{test} .

3.1.2 Supervised and Unsupervised Learning

One of the most basic separations in machine learning is the partition between supervised learning(ref) and unsupervised learning. In the case of supervised learning one knows the answer to a problem, and let the computer deduce its own logic to figure out how we get to that result. With unsupervised learning the machine is tasked with finding patterns and relationships in data sets without any prior knowledge of the system. Some authors operate with a third category, namely reinforcement learning, where the machine learns by trial-and-error. (REF)

In this thesis we only consider supervised learning. Algorithms and challenges specifically related to unsupervised learning, and reinforcement learning, is therefore not further examined.

3.1.3 Regression and Classification Problems

A respons variable can either be qualitative or quantitative in nature

3.1.4 Features and Feature Selection

Ref tilbake til "Batteries" seksjonen og utdyp hvorfor og hvorfor ikke det er gode prediktorer.

3.2 Random Forest

3.3 Support Vector Regression

Dette bør brukes og kan da skrives om.

3.3.1 Radial Basis Function Kernel

3.3.2 SVR and the Bias-Variance-Trade-off

3.4 Mean square error and Root mean square deviation

The Mean Square Error (*MSE*) can give a measure of the quality of our estimator.(ref) It is defined as

$$MSE(\epsilon) = \frac{1}{n} \sum_n^{n-1} \epsilon^2 = \frac{1}{n_{\text{samples}}} \sum_n^{n_{\text{samples}}-1} (y_i - \hat{y}_i)^2 \quad (3.1)$$

Where \hat{y}_i is the predicted value of the i -th sample, and y_i is the corresponding true value. As such it can be thought of as the average of the square of our residuals. Therefore the *MSE* can never have negative values, and smaller values mean that we have a better prediction, where at zero there is a perfect fit.

The Root mean square deviation, or root mean square Error(*RSME*), is the squared for the *MSE*:

$$RMSE = \sqrt{MSE} = \sqrt{\frac{\sum_n^{n-1} (y_i - \hat{y}_i)^2}{n}}$$

And is thus the distance, on average, of a data point from the fitted line, measured along a vertical line. The *RSME* is directly interpretable in terms of measurement units, and so is a better measure of goodness of fit than a correlation coefficient.

3.5 R^2 score - The Coefficient of Variation

In regression validation the R^2 is the standard when it comes to measuring goodness of fit. (REF=?coef Needs new ref.) In straight terms it is the proportion of the variance in the dependent variable that is predictable from the independent variable (S).

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = 1 - \frac{\sum (y_i - f_i)^2}{\sum (y_i - \bar{y}_i)^2} \quad (3.2)$$

Where y_i are the indexed response variables (data to be fitted) and f_i the predictor variables from the model with $\epsilon_i = y_i - f_i$. The average of the response variables is denoted \bar{y}_i . The second term can also be considered as the ratio of MSE to the variance (the $1/n$ factors null each other out in a fraction), or the total sum of squares (SS_{tot}).

If the residual sum of squares (SS_{res}) is low the fit is good. However, this should be compared to the spread of the response variables. After all, if the response variables are all nicely distributed close to the mean, then getting a good SS_{res} is not suspicious. We therefore do a normalization in the fraction, taking the scale of data into consideration. In the simplest polynomial fit, using a zero order polynomial (a constant), our model would just be a constant function of the mean. The sums being equal, returning unity on the fraction and the total R^2 score would be zero. In the other extreme, if the model fits perfectly, then SS_{res} would be zero and the R^2 score would be one. In this sense we have a span of possible R^2 scores between zero and one, from the baseline of the simplest model at zero, and a perfect fit at one. In contradiction to most scores the value can be negative, because the model can get arbitrarily worse, thus giving negative values. The R^2 score is useful as a measure of how good our model is at predicting future samples.

3.6 Principle Component Analysis

Principle Component Analysis (PCA) is a procedure that uses orthogonal linear transformation to reduce the amount of feature subspaces. It goes under different names in different fields, but the most recognizable might be Singular Value Decomposition. This is done by converting a set of possible correlated variables into a set of uncorrelated variables, called principle components (PC).

The PC are arranged so that the first PC has the largest variance, meaning that it accounts for as much of the variability in the data given as possible. The next PC does the same, it accounts for as much of the variability as possible with the constraint that it is orthogonal to all the former components. These orthogonal vectors are linear combinations being an uncorrelated orthogonal basis set. Graphically the shortest vectors effects the predictions the least. PCA is sensitive to the relative scaling of the original variables, so in *sklearn.decomposition.PCA* the input data is centered but not scaled for each feature before the SVD of the data is applied.

3.7 K-fold cross validation

K-folding is a cross validation technique that allows us to generalize the trends in our data set to an independent data set. In this way we can circumvent typical problems like over-fitting and selection bias.(ref=**cross-valid**)The approach for the technique is simple. Instead of doing a regression on the entire data set, it is first segmented into k number of subsets of equal size (making sure to pick out the variables randomly before distributing them to the subsets).

Now one subset can be chosen to be the 'control' or 'validation' set while the rest of the subsets are the training sets. The desirable regression is then applied on the training set, arriving at some data fitting that is the prediction. From here it is a straight forward process to analyze how well our predicted variables compare to the validation variables, for example through the R^2 score function. However, even though the subsets are picked randomly, the validation subset used could potentially not be a representative selection of the entire set. Therefore the process is repeated k times, each time using a new subset as the validation subset. After all this is done one can simply calculate the average of the scores to get the predictive power of our model. As an added benefit, since the calculations are done anyways, the average of the predictions can be used as the final fit.

Cross validation techniques are extremely useful when the gathering of new data is difficult or, sometimes, even impossible, as we are using the extra computational power at our disposal to squeeze the most amount of relevant information out of our data.

FIGURE? link to good crossvalidation. https://scikit-learn.org/stable/modules/cross_validation.html#cross-validation

3.8 Mutual Information**

Ranking of features. This can be useful.

3.9 Machine Learning × batteries

State of the art. What has been done in this field? Similar predictions of the kind we are interested in, have been done by Sandek: **sendek2017holistic**

Part III

Method

4 Some Method title

In this section we will introduce the overall approach to the research. First introducing the data set and experimental environment, before going through technicalities in the methods used to represent physical and chemical properties of the electrode materials.

All codes are written in Python 3.7 or Fortran98 and can be found here: (www.github.com/Mahscien) ((Give a real REF)). And all computations are done on a (give computer specifics.), if nothing else is noted. The tool we used, for the most part, for our data mining and data analysis is scikit-learn.(REF)

4.1 Data set and Experimental Environment

These are the void fraction, the volumetric number density and the atomic property weighted radial distribution function(AP-RDF). Other properties used where introduces in the foundation section(ref back). Lastly we quickly mentions the basics of our algorithm.

((Sendek(REF) and colleagues were right: They stated that they could not find a sufficient amount of data on ionic conductivity for a proper model to be created, or rather, they only found 40 materials that they used to make their prediction on ionic conductivity, something that our group deemed far to little for our use.))

Needed a database, found *materialsproject*! www.materialsproject.org/batteries gave a database with a good amount of information available. First of all the reduced cell formula with consistent CIF files for all voltage pairs. Secondly many different characteristics or voltage pair properties; Average voltage, Capacity, both gravimetric and volumetric, Specific Energy (Wh/kg), Energy density (Wh/l), and a measurement of the stability ($eV/atom$). Other properties that where in *materialsproject* and to some extent tested to see if they were good predictors are; space group, energy per atom, volume of the unit cell, volume change in percentage, band gap, density, total magnetization, number of sites, and elasticity.

The database contains in more then 4400 intercalation electrodes, where we have used 2291 Lithium- ion batteries, 360 Magnesium-ion batteries, 321 Natrium-ion batteries, and 481 Calcium-ion batteries for our analysis. With new compounds being added to the database continuously, including many new structural predictors. Our method were tested on these different types of intercalation batteries, both with respect to size of database and size of unit cell.

It is important to note that we have a minimum of two predictor per property of the material

at any given run. This is due to how we defined each battery to have at least one charged and one discharged state. For any given property, i.e. Volumetric number density, we have one value calculated for the charged material, and one for the discharged material. This means that we predict for a specific charged or discharged half cell configuration.

4.2 Volumetric number density

Volumetric number density, n , is used to describe concentration of countable objects. And is defined as:

$$n = \frac{\text{\#of atoms}}{\text{Volume}} \quad (4.1)$$

Where *Volume* is the volume of the unit cell.

Technically, in the volumetric number density, there is a predictor for each individual element. That is; if the material is $\text{Mg}(\text{TiO}_2)_2$ then the the number density for magnesium, titanium, and oxygen, related only to that material, will be predictors.

It is probable that such a direct measurement of a geometrical aspect would be a good predictor due to the amount of physical information. If RF were to be applied on to the entirety of the CIF file, it is probable that it would be a bad fit, due to the Bias-Variance-trade-off as mentioned in (REF), and because of the complexity of some of these files.

4.3 Void Fraction

Void Fraction, or the porosity, is a measurement of the void space in the material. Calculated *ab initio* with Poreblazer^{??}. We measure the accessible void, that is, the total amount of void space accessible from the surface. The pore volume is obtainable if Gurvich rule^{??} is fulfilled. It states that "if the density of the saturated nitrogen in the pores is assumed equal to its liquid density, regardless of the shape of the internal void network and, because of the weak interactions, regardless of the chemistry of the framework." The pore volume (v_{pore}) and the porosity (θ) are computed from:

$$v_{pore} = \frac{n_{N_2}^{ads,satd}}{\rho_{N_2}^{liq}} \quad (4.2)$$

$$\theta = v_{pore} \cdot \rho_{cryst} \quad (4.3)$$

Where $n_{N_2}^{ads,satd}$ is the specific amount of nitrogen adsorbed, $\rho_{N_2}^{liq}$ is the density of liquid nitrogen, and ρ_{cryst} is the density of the crystal in question.

Two different pore volumes are calculated, the geometric pore volume, G_{epv} , which is defined as all the free volume of the unit cell, and Helium pore volume, H_{epv} , where the unit cell that

can fit a probe with realistic intermolecular potential is tested. The calculation are done on the fixed 0K temperature.

Void Fraction is a good characterization method for microporous crystals and have had great success in metal organic frameworks (MOFS), as demonstrated also by the team of supervisors.

In case of dens materials like the one we consider in this work, the void fraction should not be a good predictor. However we decided to include it in our tests in case the space occupied by the ion in the discharged material would impact our prediction, as will be discussed later. (REF)

4.4 AP-RDF Descriptors of Electrode materials

Atomic property weighted radial distribution function (AP-RDF)?? was found to be a good predictor which also, when tested by the PCA(REF to theory part), exhibited good discrimination of geometrical and other properties, in one of their cases, gas uptake.

One of the methods found, that seemed to yield good predictions dependent on chemical properties where the Atomic property weighted radial distribution function, successfully used on MOFS. ?? Due to it looking reasonable we decided to try it out.

The radial Distribution Function(RDF) is the interatomic separation histogram representing the weighted probability of finding a pair of atoms separated by a given distance.(REF) In a crystalline solid, the RDF plot has an infinite number of sharp peaks where the separation and height are characteristic of the lattice structure. We used the minimum image convention (boundary condition)Do I need a ref here? and the RDF scores will be uniquely defined inside of the unit cell, per material-ID. The RDF can be expressed as:

$$RDF^P(R) = f \sum_{i,j}^{\text{all atom pairs}} P_i P_j e^{-B(r_{ij}-R)^2} \quad (4.4)$$

In our case the RDF scores in a electrode framework has been interpreted as the weighted probability distribution to find a atom pair in a spherical volume of radius R inside the unit cell according to equation above. 4.4

Summing over all the atom pairs, where R_{ij} is the minimum image convention distance of these pairs, B is a smoothing parameter, and F is a scaling or normalization factor. Our Own approach to this is written in Fortran, and can be found in the appendix with an operational pdf.(REF)

The RDF can be weighted to fit the requirements of the chemical information to be represented, by introducing the atomic properties, P_i and P_j . We weighted the radial probabilities by three tabulated atomic properties namely electronegativity, polarizability, and Van der Waals volume, which gives us the AP-RDF. While a regular RDF function encodes geometric features, the atomic property weighted RDF additionally characterizes the chemical features within a material. An atomic property weighted RDF can be seen on the screen.

To test our method, we used it to reproduce the results for the two MOFS, namely *IRMOF-1* and *MIL-45* found in the article by Fernandez.?? We confirmed their findings .. though with drawback related to the size... which are flawed in our case. In our opinion, we think that this is a fundamental drawback, and the results depends on the size of the simulation cell(which can be made by replicating the unit cell).

INSERT BILDET AV PLOT AV AP-RDF.

4.5 Algorithm

First of all, a wrote a program for "scraping" the *materialsproject* webpage for batteries(0). This gave us the possibility to gather all the available resources on the batteries in the database in a fast and effective manor, as well as updating these CSV files of battery-IDs. (ref)

We then run a second program that downloads all the information on the materials that matches a material-ID correlated to a battery-ID(1,2). (ref) Before constructing a CIF file structured so that all the battery-IDs, charged-material-IDs, and dischargerd-material-IDs are correlated with the information on the charged and discharged properties.

After, the volumetric density fraction is calculated(3) and added to the main CSV file for both charged and discharged materials. While the CIF files are being processed for Pore-blazer(4) where the void fraction is calculated(5,6).

Then we merge all our CSV files based on what properties that we are interested in and makes a CSV file called *for_ML.csv*(7,8) that we feed into our random forest algorithm(9). We then run cross validation, MSE, and plot what we are interested in(10).

In addition we also tested for different machine learning algorithms, as mentioned(ref), but these were only to test the reliability of our model, and will be discussed in the discussion section(ref)

maybe add this in the appendix?

Algorithm:

Steps for use of python scripts:

`mp_battery_scraper.py`

0: Scrape batteries with a given working ion from the Materials Project battery explorer (<https://www.materialsproject.org/#search/batteries>)

`fillproperties.py`

1: Download all materials that match a material_id correlated to a battid.
Output files: directory `cif_info_dir/<material_id>_prop.dat`

`add_features.py`

2: Gets and adds the material specific features from the JSON dump to a csv.
Output files: `material_properties.csv`

`elements.py`

3: Calculate the density fractions for all materials.
Output files: `out_csv_dis.csv`

`forPoreblazer.py`

4: Download the CIF files as JSON for all materials correlated to a battid.
Output files: directory `cif_for_poreblazer/<material_id>_cif.dat`

`process_cif.py`

5: Extract the CIF information from the previous JSON data.
Output files: directory `cif_for_poreblazer/cif_files/<material_id>_cif.dat.csv`

`process_cif.py`

6: Extract void fraction with poreblazer using the CIF files.
Output files: `helvol_geomvol.csv`

`merger.py`

7: Merge charged and discharged for all properties
Output files: `allFiles.csv`

`prep_csv.py`

8: Select predictors and targets for ML

Output files : for_ML.csv

randomforest.py

9: Run randomforrest

Output files : Depending on what being saved: ./Results/*

crossvalidation.py

10: Run cross-validation , remove outliers .

11: ???

12: Profit!

<i>Target: Accuracy:</i>	Average Voltage	Gravimetric Capacity	Volumetric Capacity	Specific Energy
R^2 -score	-0.0461	0.3426	0.3784	-0.0011
R^2 -train	0.8676	0.8799	0.9052	0.8764
CV:	-0.7976(+/- 1.3547)	0.1844 (+/- 0.2182)	0.2983 (+/- 0.3636)	-0.3742 (+/- 0.9904)
MSE:	1.1457	4521.6902	62785.2	92593.3
CV-mean:	-0.5899	0.1660	0.3166	-0.2962

Figure 1: Wall of numbers.

Part IV

Result & Discussion

5 Result section title

MSE, PCA, R2, compair. Ordered by target:

5.1 Random factors from database.

5.1.1 Average Voltage

5.1.2 Capacity

5.1.3 Energy Density

5.2 Volumetric number density

5.2.1 Average Voltage

Charged:

Discharged:

5.2.2 Capacity

5.2.3 Energy Density

5.3 Void fraction

5.4 AP-RDF

PCA, R2, MSE.

5.5 Stability

This is a novel work, the aim is therefor to explore different predictors, by figuring out the different weights of the predictors on different targets, and which predictors that does not favorable for our predictions.

The model is decent at predicting; Gravimetric and Volumetric Capacity(87%), Specific Energy(70%), and Energy density(68%), but has no capability of predicting stability as of now. There is a need for *ab initio* calculations for several of our predictors, they calculate something that we know most definitely is correlated to the target without any premature calculations. This is something we are trying to move away from. As of now, only using the density fraction, we can get somewhere between 40% to 60% accuracy with our model.

This did not improve when including the void fraction, our predicitions actually got worse.

AP-RDF - Still no good results! **This is a struggle.**

5.6 Geometrical descriptors

These grafs all represent the accuracy of the predictions on the training data and on new data given to the machine, with only the number density as a predictor, and the Average voltage, Gravimetric capacity, Volumetric capacity, energy density, and physical stability for the discharged material, as targets (a-f)???. Most notably the predictions on the Average voltage, Gravimetric capacity, Volumetric capacity, energy density, and specific energy are all showing a decent amount of correlation, with around 60% accuracy. The physical stability for the discharged-,and the physical stability of the charged-materials show that there is no correlation between the number density and the physical stability. It is also shown that there is no correlation between the number density and the void fraction. Or any of the other properties for that matter.

Part V

Summary

6 Summary and future work

6.1 Batteries

6.2 future work

6.2.1 improving method