

PREDICTING
SOLID-STATE QUBIT
MATERIAL HOST

by

Oliver Lerstøl Hebnes

THESIS
for the degree of
MASTER OF SCIENCE



Faculty of Mathematics and Natural Sciences
University of Oslo

April 25, 2021

Contents

I	Introduction	1
1	Introduction	3
1.1	Motivation	4
1.2	Holy grail	4
1.3	Structure of thesis	4
II	Theory	5
2	Quantum technologies	7
2.1	Quantum computing	8
2.1.1	Quantum computing requirements	10
2.2	Quantum communication	10
2.3	Quantum sensing	11
2.4	Available quantum platforms	12
2.5	Introduction to semiconductor physics	14
2.5.1	Point defects in semiconductors	17
2.5.2	Optical defect transitions	18
2.6	Semiconductor candidates for quantum technology	20
2.6.1	Diamond - the benchmark material for QT	20
2.6.2	Qubit material host requirements	23
2.6.3	Silicon carbide	25
2.6.4	Alternative promising material hosts	26
2.6.5	Associated challenges with material host discovery	27
3	Understanding information from ab-initio calculations	29
3.1	Introduction to density functional theory	29
3.1.1	The Schrödinger equation	30
3.1.2	The many-particle Schrödinger equation	30
3.1.3	The Born-Oppenheimer approximation	32
3.1.4	The Hartree and Hartree-Fock approximation	32
3.1.5	The variational principle	33
3.2	The density functional theory	33

3.2.1	The Hohenberg-Kohn theorems	34
3.2.2	The Kohn-Sham equation	34
3.2.3	The exchange-correlation energy	36
3.2.4	Limitations of the DFT	38
4	Machine learning	41
4.1	Supervised learning	42
4.2	Evaluating accuracy of a model	42
4.2.1	Bias-variance tradeoff	43
4.2.2	Accuracy, precision and recall	44
4.2.3	Cross validation	47
4.3	Logistic regression	48
4.3.1	Stochastic gradient descent	49
4.4	Decision trees	50
4.4.1	Growing a classification tree	51
4.4.2	Classification algorithm	52
4.4.3	Pruning a tree	52
4.4.4	Pros and cons of decision trees	52
4.5	Ensemble methods	53
4.5.1	Bagging	53
4.5.2	Boosting	54
4.6	Dimensionality reduction	57
4.6.1	Principal component analysis	57
4.7	Practical challenges associated with machine learning	60
III	Methodology and implementation	63
5	Material Science Databases	65
5.1	Fundamentals of a database	65
5.1.1	API and HTTP requests	65
5.1.2	Open-source Python libraries for material analysis	67
5.2	Databases and cloud services	69
5.2.1	Novel Materials Discovery	69
5.2.2	Materials project	69
5.2.3	AFLOW	70
5.2.4	Open Quantum Materials Database	70
5.2.5	JARVIS	71
6	Information flow	73
6.1	Extraction and featurization of data	73
6.1.1	Practical data extraction with Python-examples	74
6.2	Matminer featurization	80

6.3	Labelling of data	83
6.3.1	First approach; the Ferrenti approach	83
6.3.2	Second approach; the augmented Ferrenti approach	86
6.3.3	Third approach; the insightful approach	87
6.3.4	Comparison of the approaches	90
6.4	Model selection	92
IV	Results and discussion	95
7	Optimization	97
7.1	Comparing functionals for band gaps	99
7.2	Technical details on ML classifiers	103
7.2.1	The Ferrenti approach	103
7.2.2	The augmented Ferrenti approach	106
7.2.3	The insightful approach	108
V	Appendices	111
A	Extracting information from ab-initio calculations	113
A.1	The Born-Oppenheimer approximation	113
A.2	The variational principle	114
A.3	The Hohenberg-Kohn theorems	115
A.3.1	The Hohenberg-Kohn theorem 1	115
A.3.2	The Hohenberg-Kohn theorem 2	117
A.4	Self-consistent field methods	117
B	Featurization	119
B.1	Table of featurizers	119
B.2	Erroneous entries	124
B.3	Optimization	125

Part I

Introduction

Chapter 1

Introduction

Discovery of current deep defects and material hots by serendipitet. (ref. data-mining artikkl)

This it The introduction. Another coffee.

Trengs det egentlig undertitler her? Tenketenketenk

1.1 Motivation

1.2 Holy grail

1.3 Structure of thesis

Part II

Theory

Chapter 2

Quantum technologies

This chapter will provide a brief overview of the current state-of-the-art in quantum technological advances. This will not only give us insights in how the technology is being used today, but also grant us the opportunity to discuss key concepts that are fundamental to understand for this thesis. Thereafter we will look into how materials are built up, and what kind of properties a material needs to exhibit to be an eligible host for quantum devices. Finally, we will give a few specific examples of materials with promising point defects that have been comprehensively researched. Importantly, this will motivate the reasoning for finding new materials that might excel in areas where other materials fall short for utilization in quantum technology.

Quantum technology (QT) refers to practical applications and devices that utilize the principles of quantum physics as a foundation. Technologies in this spectrum are based on concepts such as *superposition*, *entanglement* and *coherence*, which are all closely related to one another.

A quantum superposition refers to that any two or more quantum eigenstates can be added together into another valid quantum state, such that every quantum state can be represented as a sum, or a superposition, of two or more distinct states. This is according to the wave-particle duality which states that every particle or another quantum entity may be described as either a particle or a wave. When measuring the state of a system residing in a superposition of eigenstates, however, the system falls back to one of the basis states that formed the superposition, destroying the original configuration.

Quantum entanglement refers to when a two- or many-particle state cannot be expressed independently of the state of the other particles, even when the particles are separated by a significant distance. As a result, the many-particle state is termed an entangled state [1].

Quantum coherence arises if two waves coherently interfere with each other and generate a superposition of the two states with a phase relation. Likewise, loss of coherence is known as *decoherence*.

Another concept that the reader should be familiar with is the famous Heisenberg uncertainty principle. It states that

$$\sigma_x \sigma_p \leq \frac{\hbar}{2}, \quad (2.1)$$

where σ_x is the standard deviation for the position and σ_p is the standard deviation in momentum. This means that we cannot accurately predict both the position and momentum of a particle at the same time. Thus, we often calculate the probability for a particle to be in a state which results in concepts such as an electron sky surrounding an atom core. However, remember that equation (2.1) is an inequality, which means that it is possible to create a state where neither the position nor the momentum is well defined.

2.1 Quantum computing

The start of the digital world's computational powers can be credited to Alan Turing. In 1937, Turing [2] published a paper where he described the *Turing machine*, which is regarded as the foundation of computation and computer science. It states that only the simplest form of calculus, such as boolean Algebra (1 for true and 0 for false), is actually computable. This required developing hardware that could handle classical logic operations, and was the basis of transistors that are either in the state ON or OFF depending on the electrical signal. Equipped with a circuit consisting of wires and transistors, commonly known as a computer, we could develop software to solve all kinds of possible applications.

Driven by the development of software, conventional computers have in accordance to Moore's law [3], doubled the amount of transistors on integrated circuit chips every two years as a result of smaller transistors. Furthermore, the clock frequency has enhanced with time, resulting in a doubling of computer performance every 18 months [4]. Alas, miniaturization cannot go on forever as transistors are mass-produced at 5 nm today and are expected to reach a critical limit of 3 nm in the following years [5].

To sustain the digital world's increasing computational demand, other alternatives than the conventional classical computer must be explored. This is where quantum computing comes into the picture. The term quantum computer is a device that exploits quantum properties to solve certain computational problems more efficiently than allowed by Boolean logic [6].

The idea is to pass information in the form of a quantum bit, or *qubit* for short. They are the building blocks of quantum computers, and as opposed to the conventional 0 or 1-bits that classical computers are based on, they can inhabit any superposition of the states 0 or 1. This is illustrated in figure 2.1.

The architecture of a gate-based quantum computer is dependent on a set of quantum logic gates that perform unitary transformations on sets of qubits [7, 8]. Other implementations of quantum computers exist, such as the adiabatic quantum computer. This approach is not based on gates, but on defining the answer of a problem as the ground state of a complex network of interactions between qubits, and then controlling the interactions to adiabatically evolve the system to the ground state [9].

It has been demonstrated that exponentially complex problems can be reduced to polynomially complex problems for quantum computers [4]. For example, a quantum search algorithm found by Grover [10] offers a quadratic speed-up compared to classical algorithms, while Shor's quantum integer factorization algorithm [11] presents an exponential speed-up. Intriguingly, Google reported in 2019 that they ran a random number generator algorithm on a superconducting processor containing 53 qubits in 200 seconds, which would most likely take several times longer for a classical supercomputer to solve [12]. It is anticipated that quantum computers will excel in exceedingly complex problems, while many simpler tasks may not see any speed-up at all compared to the classical regime. Hence, quantum- and classical computers are envisioned to coexist for each their purpose.

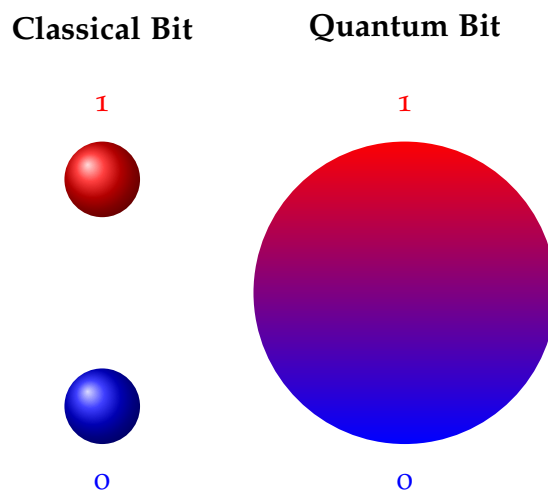


Figure 2.1: Conceptual illustration of the two-level classical bit, which are restricted to the boolean states 1 (true) or 0 (false), and the quantum bit that can be in any superposition of the states 0 or 1.

Quantum computing is a highly sought-after goal, but there are extensive challenges that need to be addressed. Controlling a complex many-qubit system is difficult, since it is not always possible to establish interactions between qubits [7] and maintain entanglement over both time and distance. Additionally, decoherence and other quantum noise occurs as a result of the high volatility of quantum states, making quantum state manipulation prone to errors. The *quantum error correction* protocols and the theory of *threshold theorem* deals with this vulnerability, stating that noise most likely does not pose any fundamental barrier to the performance of large-scale computations [4].

2.1.1 Quantum computing requirements

As ever-promising the concepts of quantum technology are, the physical realizations are in the preliminary stage of development. Here we will concretize critical principles for a physical realisation of a quantum platform.

“I always said that in some sense, these criteria are exactly the ones that you would teach to kindergarten children about computers, quantum or otherwise” DiVincenzo [13]

DiVincenzo formulated in the year of 2000 seven basic criteria for a physical qubit system with a logic-based architecture [7].

1. A scalable physical system with well characterized qubits
2. The ability to initialize the state of the qubits to a simple initial system
3. Have coherence times that are much longer than the gate operation time
4. Have a universal set of quantum gates
5. Have the ability to perform qubit-specific measurements

These five criteria must be met for a quantum platform to be considered a quantum computer.

2.2 Quantum communication

Quantum communication refers to the transfer of a state of one quantum system to another. Since information can be stored in qubits, we picture *flying qubits* that transfer information from one location to another [14]. The benefits of using flying qubits are in particular valued in quantum cryptography, since the quantum nature of qubits can be exploited to add extra layers of security [4].

Consider the example of encrypting a digitally transmitted conversation. It is difficult to avoid someone eavesdropping on a conversation, however, the problem is diminished if the eavesdropper does not speak the language, keeping the information in the conversation safe. This is the original idea of encryption, such that the information has been encrypted into something incomprehensible for any eavesdropper. A common practice is to encrypt information and share a public key, which everyone can read, and a private key, only known for the sender and receiver of information. This should be sufficient to keep the information secure, given that the complexity of the private key is impenetrable.

Importantly, we live in a digital world where most of our actions are increasingly being stored as information, and we could imagine that the eavesdropper in the latter example stored the conversation. Even if the content of the conversation was encrypted, it still presents a challenge, since encrypted information stored today could be deciphered in ten or twenty years' time. Consequently, finding an encryption method that could make information either impossible to eavesdrop on or make the security unbreakable forever is very desirable. This is the ultimate goal of quantum cryptography [4].

Consider the example of information encoded into a qubit as a superposition of two quantum states. Now, if a wild eavesdropper would try to measure the information, the nature of quantum physics tells us that the original configuration would be destroyed and the receiver would be alerted of the eavesdropper. Furthermore, if the eavesdropper would try to make a copy of the message, the copying itself would be limited of the no-cloning theorem [15] which declare that quantum states cannot be copied.

A clever approach to ensure confidentiality is to send the encryption key before sending the actual encrypted information. If the key is received unperturbed, the key remains secret and can be safely employed. If it turns out perturbed, confidentiality is still intact since the key does not contain any information and can be discarded. This approach is termed the *quantum key distribution* (QKD) [15, 16]. It should be noted that this requires both the sender and receiver to have access to methods for sending, receiving and storing qubit states, such as a quantum computer. Additionally, the sender and receiver will need to initially exchange a common secret which is later expanded, making quantum key *expansion* a more exact term for QKD [4, 16].

Most applications and experiments use optical fibers for sending information via photons, with the distance regarded as the main limitation. This is because classical repeaters are unable to enhance quantum information because of the no-cloning theorem, making photon loss in optical fiber cables inevitable. Thus, quantum communication must reinvent the repeater concept, using hardware that preserves the quantum nature [17] and are compatible with wavelengths used in telecommunication. Nonetheless, secure QKD up to 400 km has recently been demonstrated using optical fibres in academic prototypes [18].

2.3 Quantum sensing

Measurements are part of our digital world today to a great extent. There would be no way to exchange goods, services or information without reliable and precise measurements [17]. Thus, improving the accuracy of sensors for every measurement done is desirable. One method to improve measurement accuracy, resolution and sensitivity can be by utilizing quantum sensors.

Quantum sensors exploit quantum properties to measure a physical quantity [19]. This is possible because quantum systems are highly susceptible to perturbations to its surroundings, and can be used to detect physical properties such as either temperature or an electrical or magnetic field [19].

For a quantum system to be able to function as a quantum sensors, a few criterias needs to be met. Firstly, the quantum system needs to have discrete and resolvable energy levels. The quantum system also needs to be controllably initialised into a state that can be identified and coherently manipulated by time-dependent fields. Lastly, the quantum system needs to be able to interact with the physical property one wants to measure through a coupling parameter [19].

It is also possible to also exploit quantum entanglement to improve the precision of a measurement. This gain of precision is used to reach what is called the Heisenberg-limit, which states that the precision scales as the number of particles N in an idealized quantum system [17, 19], while the best classical sensors scale with \sqrt{N} .

2.4 Available quantum platforms

Many different quantum platforms have been physically implemented, and this section will serve as a brief overview of the current status. For a more thorough review of qubit implementations, the reader is directed to Refs. [8, 17].

Superconducting circuits can be used in quantum computing, since electrons in superconducting materials can form Cooper pairs via an effective electron-electron attraction when the temperature is lower than a critical limit. Below the limit, electrons can move without resistance in the material [20]. Exploiting this intrinsic coherence, qubits can be made by forming microwave circuits based on loops of two superconducting elements separated by an insulator, also known as Josephson tunnel junctions [17]. Today, superconducting Josephson junctions are the most widely used quantum platform, but they requires very low temperature (mK) to function, making them costly to use. Additionally, the current devices experience a relatively short coherence time, causing challenges in scaling up.

Single photons is an eligible quantum platform that can be implemented as qubits with one-qubit gates being formed by rotations of the photon polarization. Its use in fiber optics are less prone to decoherence, but faces challenges since the more complex photon-photon entanglement and control of multi-qubits is strenuous [8].

By fixing the nuclear spin of solid-state systems, it is possible to implement a quantum platform that experience long spin coherence. This enables the manipulation of qubits that utilize electromagnetic fields, making one-qubit gates realizable.

The isolated atom platform is characterized by its well-defined atom isolation. Here, every qubit is based on energy levels of a trapped ion or atom. Quantum entanglement can be achieved through laser-induced spin coupling, however scaling up to large atom numbers induce problems in controlling large systems and cooling of the trapped atoms or ions.

A quantum dot (QD) can be imagined as an artificial atom which is confined in a solid-state host. As an example, a quantum dot can occur when a hole or an electron is trapped in the localized potential of a semiconductor's nanostructure. QDs exhibit similar coherence potential as the isolated atom platform, but without the drawback of confining and cooling of the given atom or ion [17]. Moreover, it is possible to limit decoherence due to nuclear spins by dynamic decoupling of nuclear spin noise and isotope purification [8].

A QD can normally be defined lithographically using metallic gates, or as self-assembled QDs where a growth process creates the potential that traps electrons or holes. The difference between them is a question of controllability and temperature, since the metallic gates is primarily controlled electrically and operate at < 1 K, while self-assembly QDs are primarily controlled optically at ~ 4 K [8]. Despite requiring very low temperatures, QDs have the potential for fast voltage control and optical initialization. As with trapped ions, electrostatically defined quantum dots experience a short-range exchange interaction, imposing a limitation for quantum computing and quantum error correction protocols. A potential solution could include photonic connections between quantum dots. On the contrary, self-assembled quantum dots couple strongly to photons due to their large size in comparison to single atoms. However, the size and shapes of self-assembled quantum dots are decided randomly during the growth process, causing an unfavourable large range of optical absorption and emission energies [8].

Lastly, we will turn towards point defects in bulk semiconductors as a physical implementation of a quantum platform. Point defects shares many of the attributes of quantum dots, such as discrete optical transitions and controllable coherent spin states, but are vulnerable to small changes in the lattice of the semiconductor. Thus, it can be difficult to isolate a point defect from the surrounding environment. However, one can utilize the strength of the solid-

state semiconductor host to isolate to some extent the point defect, yielding extended coherence times and greater optical homogeneity than other quantum dot systems. Before we dwell into the intricacies of point defect qubits as a building block for QT, we will provide the necessary background for the crystal- and electronic structure of semiconductors.

2.5 Introduction to semiconductor physics

The interactions between atoms and characteristics of matter form the foundation of materials science. The applications of materials science are extensive, with examples such as a bottle of water or to a chair to sit in.

Solid materials, like plastic bottles, are formed by densely packed atoms. These atoms can randomly occur through the material without any long-range order, which would categorize the material as an *amorphous solid*. Amorphous solids are frequently used in gels, glass and polymers [21]. However, the atoms can also be periodically ordered in small regions of the material, classifying the material as a *polycrystalline solid*. All ceramics are polycrystalline with a broad specter of applications ranging from kitchen-porcelain to orthopedical bio-implants [22]. A third option is to have these atoms arranged with infinite periodicity, making the material a *crystalline solid* or more commonly named a *crystal*. The three options are visualised in figure 2.2. Hereon, we will focus on crystalline solids.

The periodicity in a crystal is defined in terms of a symmetric array of points in space called the *lattice*, which can be simplified as either a one-dimensional array, a two-dimensional matrix or a three dimensional vector space, depending on the material. At each lattice point we can add an atom to make an arrangement called a *basis*. The basis can be one atom or a cluster of atoms having the same spatial arrangement. Every crystal has periodically repeated building blocks called *cells* representing the entire crystal. The smallest cell possible is called a *primitive cell*, but such a cell only allows lattice points at its corners and it is often quite rigid to work with when the structure becomes complex. As a solution, we will consider the *unit cell*, which allows lattice points on face centers and body centers.

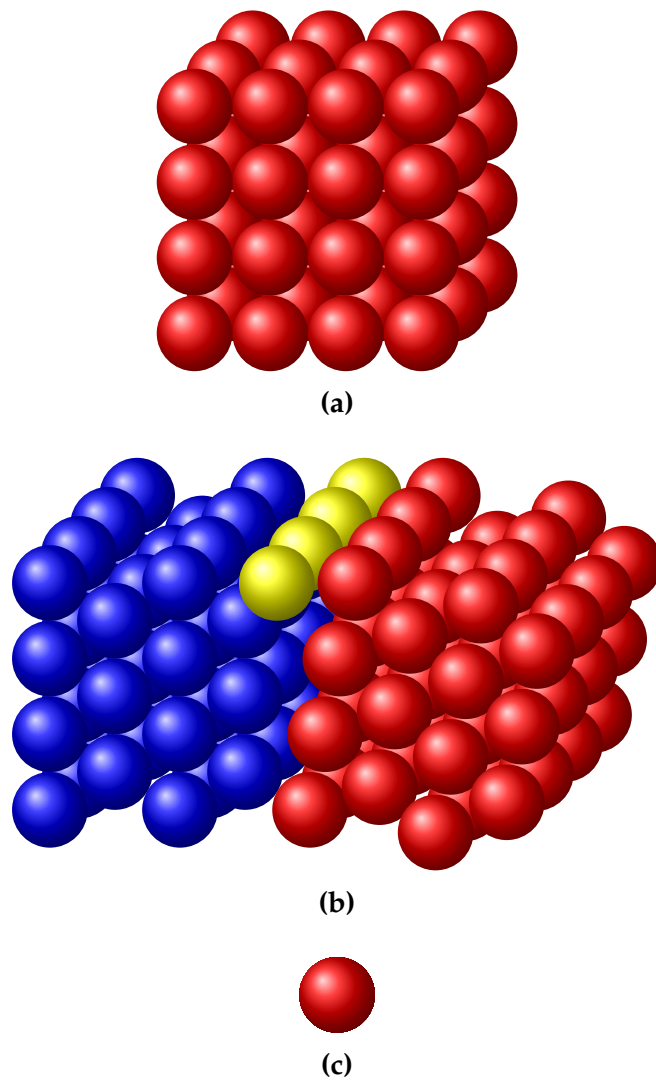


Figure 2.2: Schematic representation of different degrees of ordered structures, where (a) is a crystalline of a simple cubic lattice, (b) is a polycrystalline hexagonal lattice, and (c) is an amorphous lattice.

One example of a crystal structure is the perovskite structure. Compounds with this structure are characterized by having an ABX_3 stoichiometry whose symmetries belong to one of 15 space groups identified by Lufaso & Woodward [23], such as the cubic, orthorhombic and tetragonal. For our purpose, we will be looking into when the X atom is oxygen, and refer to the oxygen-perovskite ABO_3 . The A atom is nine- to 12-fold coordinated by oxygen, while the B atom is sixfold coordinated by oxygen, and the BO_6 octahedra are connected to the corners in all three directions as visualized in figure 2.3.

The motivation behind the research on perovskites is related to the large amount of available ABO_3 chemistries, where a significant portion of these take the perovskite structure. Perovskites have a broad specter of applications, ranging from high-temperature superconductors [24] and ionic conductors [25] to multiferroic materials [26]. Additionally, adding a perovskite-type compound to solar cells has reportedly resulted in higher performance efficiencies while being cheap to produce and simple to manufacture [27, 28]. However, this includes the use of hybrid organic-inorganic compounds and excludes the use of oxygen.

Isolated atoms have distinct energy levels, where the Pauli exclusion principle [29] states for fermions that each energy level can at most accommodate two electrons of opposite spin. In a solid, the discrete energy levels of the isolated atom spread into continuous energy bands since the wavefunctions of the electrons in the neighboring atoms overlap. Hence, an electron is not necessarily localized at a particular atom anymore. This is exemplified as every material has a unique band structure, similar to every human having their unique fingerprint.

Knowing which energy bands are occupied by electrons is the key in understanding the electrical properties of solids. The highest occupied electron band at 0 K is called the valence band (VB), while the lowest unoccupied electron band is called the conduction band (CB). The energy gap of forbidden energy levels between the maximum VB and the minimum CB is known as the band gap, and its energy is denoted as E_g . If a material can be classified

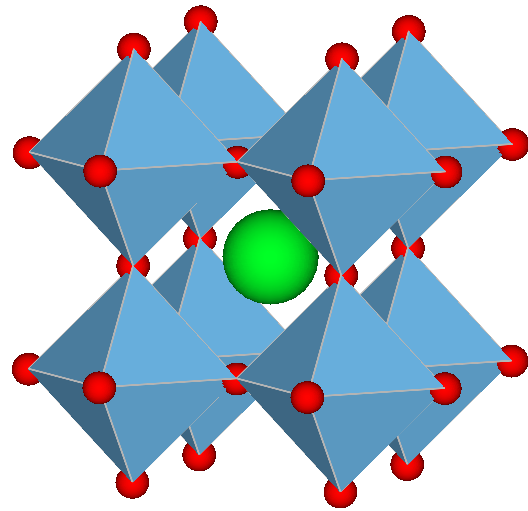


Figure 2.3: A crystal structure of $SrTiO_3$ which is a cubic perovskite. The red atoms are oxygen, whereas the green atom is strontium, and inside every corner-sharing BO_6 octahedral unit is a titanium atom.

as a semiconductor depends on the band gap and the electrical conductivity. As an example, Silicon is commonly thought of as a semiconductor, and has a band gap of about 1.12 eV at 275 K [30].

To be able to accelerate electrons in a solid using an electrical field, they must be able to move into new energy states. At 0 K, the entire valence band of a semiconductor is full with electrons and there are no available states nearby, making it impossible for current to flow through the material. This can be solved by using either thermal or optical energy to excite electrons from the valence band to the conduction band, in order to *conduct* electricity. At room temperature, some semiconductors will have electrons excited to the conduction band solely from thermal energy matching the energy band gap [21].

In some scenarios, thermal or optical energy is not sufficient for an excitation since the energy bands are also dependent on the crystal momentum. A difference in the momentum of the minimal-energy state in the conduction band and the maximum-energy state in the valence band results in an *indirect bandgap* as seen in figure 2.4a. If there is no difference at all, the material has a *direct bandgap*, which is visualized in figure 2.4b.

Electrons in semiconductor materials can be described according to the Fermi-Dirac distribution

$$f(E) = \frac{1}{1 + e^{(E-E_F)/kT}},$$

where k is Boltzmann's constant, T is temperature, E is the energy and E_F is the Fermi level. The Fermi-Dirac distribution gives the probability that a state will be occupied by an electron, and at $T = 0$ K, every energy state lower than E_F is occupied by electrons while the opposite is true for energy states above E_F [21].

2.5.1 Point defects in semiconductors

In real life, a perfect crystal without any symmetry-breaking flaw does not exist. These flaws are known as defects and can occur up to three dimen-

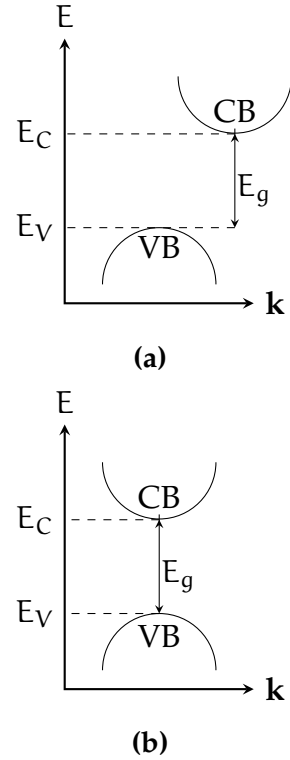


Figure 2.4: A schematic drawing of (a) an indirect- and (b) a direct bandgap.

sions. An example one-dimensional defect is known as a *line defect*, while two dimensional defects can be *planar defects*, and in three dimensions we have *volume defects*. Lastly, defects can also occur in zero dimensions and are then termed *point defects*. Point defects normally occur as either vacancies, interstitial placement inbetween lattice sites or as substitution of another existing atom in the lattice.

Defects can greatly influence both the electronic and optical properties of a material. A substitutional defect can at first be regarded as an impurity or an antisite, but they can also be intentionally inserted, an approach known as *doping*. Doping can result in an excess of electrons or holes, making the semiconductor either an n- or p-type, respectively. Consequently, the semiconductor will have energy levels in the (forbidden) band gap that originates from the defects. If the energy levels introduced are closer than ~ 0.2 eV to the band edges, they are termed *shallow* defects.

Shallow defects can contribute with either excess electrons to the conduction band, or excess holes to the valence band. However, the induced charge carriers (electrons or holes) interact strongly with the band edges, resulting in a delocalized wavefunction regarding the position in the lattice.

For the opposite case, if the energy levels rests closer to the middle of the semiconductor's gap, the introduced defects are known as *deep level* defects. Deep levels normally occur due to either dangling bonds or impurities, and have highly localized electron wavefunctions. This might assure the isolation required for long coherence times, which is an appealing promise in quantum technological advances.

Deep levels can be unfortunate in semiconductors since they can interact with the charge carriers, potentially destroying the desired electronic or optical property of the material. Deep level defects can function as electron-hole recombination centers, or to trap charge carriers, yielding the commonly used name deep level *traps*. Both of the given situations results in a lower concentration of charge carriers, which showcase why deep levels can be unwanted in semiconductor devices. However, deep level defects show extraordinary properties in quantum technology due to their ability to ensure isolation and preserve coherence.

2.5.2 Optical defect transitions

Optical transitions refers to excitation of charge carriers due to either emission or absorption of electromagnetic radiation, and can be done with a laser light or electron beam. Figure 2.5 represents a configuration coordinate (CC) diagram of a defect transition. The y-axis is a function of the energy E , while the x-axis is a function of the configuration coordination Q . The lowest point in the lower parabola is known as the ground state (GS) configuration Q_{GS} ,

which is the most stable atomic position, while for the upper parabola it is known as the excited state configuration Q_{ES} . The dotted lines represent vibronic excitations to the energy of the ground state Q_{GS} for the lower parabola, while it represents Q_{ES} for the higher parabola.

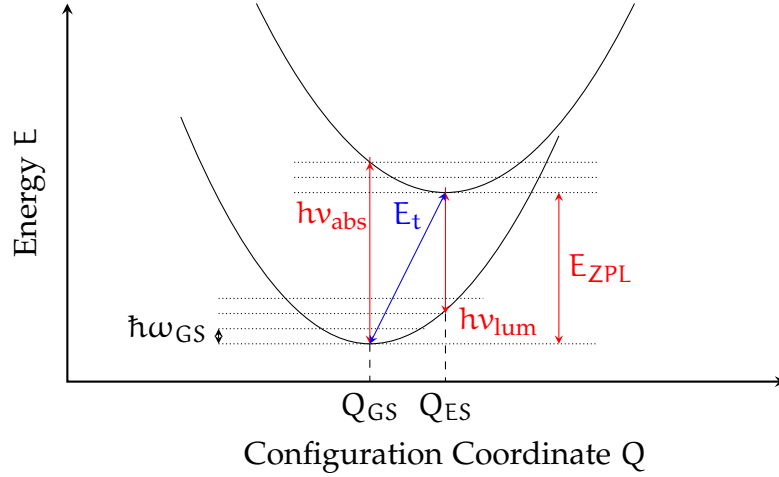


Figure 2.5: A schematic representation of a configuration coordination diagram based on Ref. [31].

The optical transitions in figure 2.5 are marked with red arrows. During slow transitions, such as during thermodynamic defect transitions, the original configuration have time to rearrange due to phonon vibrations. This is schematically drawn as the blue arrow, where the energy E_t equals the ionization energy or the position of the defect level. Optical transitions, on the other hand, are marked in red and occur in a short time range such that the original configuration does not change. They can appear in the exchange of charge carriers with the band edges, and in a defect's internal excited state, with the latter scenario being most relevant for this thesis.

Consider a defect that rests in the ground state configuration Q_{GS} . Suddenly, it absorbs a photon with energy $h\nu_{abs}$ and occupies an excited vibronic state of the upper parabola after a vertical transition. Through lattice reconfigurations, the defect will move towards the bottom of the upper parabola, also known as Q_{ES} . Eventually, it will relax to the lower parabola by emitting a photon with energy $h\nu_{lum}$, also known as a zero-phonon line (ZPL) of energy E_{ZPL} . On the other hand, any transitions between vibronic excitation levels are phonon-related. How strong the electron-phonon interaction is can be quantified by the Huang-Rhys factor S [32]. If the two parabolas in figure 2.5 have the same configuration of Q , emission into the ZPL is enabled and $S \sim 0$. The stronger the coupling, the smaller amount of emission in the ZPL.

The optical properties of a host material can be greatly influenced by defects, in particular the ES to GS transition that can occur in a defect, as dis-

cussed for figure 2.5. If the defect were to facilitate the emission of single photons with a detectable time inbetween together with a distinguishable ZPL, the defect would be referred to as a single photon source (SPS). The criteria for SPS are not met in many materials, since charge-state transitions often comprise interactions with either the VB or the CB. Thus, most SPSs' GS and ES levels are situated within the band gap of a host material. Consequently, mostly wide-band gap semiconductors are used as host materials for SPSs.

2.6 Semiconductor candidates for quantum technology

The properties of point defects are promising in a quantum technological perspective. We have seen that point defects facilitate deep energy levels within the band gap of the semiconductor, and provide isolation in the solid-state matrix as a result from a high degree of localization of the defect orbitals. If the host material have a small spin-orbit coupling, it could provide long coherence times for a deep level trap in localized and high-spin states. Additionally, point defects have the potential to be single-photon sources, giving rise to sharp and distinguishable optical transitions, where a significant amount of the emission can be of the energy E_{ZPL} . This is in particular seen in wide-bandgap semiconductors, and combined with a weak electron-phonon interaction, can have the capacity to be fabricated as a high-fidelity SPS with a significant ZPL part.

In this section we will provide specific examples of a variety of promising candidates, and what properties they possess that makes them auspicious. Additionally, we will briefly mention what the challenges with the candidates are, and why it is important to explore other viable options.

2.6.1 Diamond - the benchmark material for QT

The most studied point defect system is the nitrogen-vacancy (NV^{-1}) in diamond. Figure 2.6 schematically shows the different stages of constructing the negative charge state. Panel 2.6a shows the electronic states that correspond to the difference for an isolated atom and a lattice of atoms, as a superposition of sp^3 orbitals that generates valence and conduction bands. In panel 2.6b, a vacancy has been created by removing a carbon atom, and the four orbitals interact with each other resulting in two new states with a_1 and t_2 symmetry due to dangling bonds. Substituting a carbon atom with a nitrogen atom further splits the t_2 -states into two new states. The states $a(1)$ and e_x, e_y are of importance, as they are the GS and the ES of the qubit defects,

respectively. Here, an optical spin-conserving transition can occur due to a laser light of correct wavelength [33], as exemplified from the discussion from the last section.

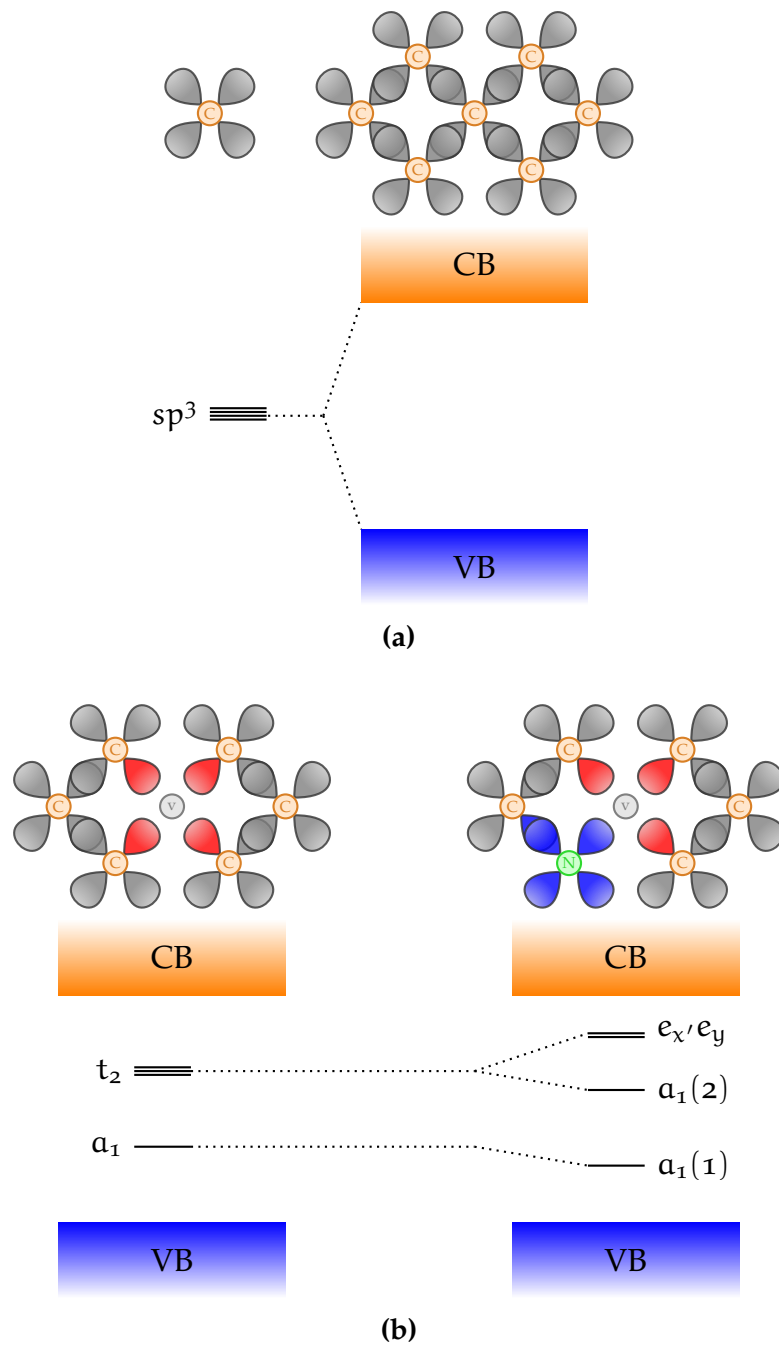


Figure 2.6: A schematic representation of the electronic structure of the NV⁻¹ defect in a tetrahedrally coordinated semiconductor, exemplified by diamond. Figure adapted from Ref. [33].

The nitrogen-vacancy in diamond is a prominent single-photon source up to room temperatures. This involves initializing, manipulating and reading out of the qubit state using optical and electric excitations, and electric and magnetic fields [33]. The potential qubit system have promising applications in quantum- communication and computation, with a demonstrated entanglement between two NV center spins that are separated by 3 m [34]. Nevertheless, perhaps the most propitious application can be seen in quantum sensing as high-sensitivity magnetometer with nanoscale resolution [35].

Unfortunately, the NV-center display restricted capabilities for quantum communication and computation. The amount of emission into the zero-phonon line is 4% at 6 K [36], which is low. The emission of the qubit center is not completely compatible with current optical fiber technologies, since the emission is in the red wave-length specter. Additionally, fabricating materials of diamond is far from unchallenging and serves as a significant incentive to find other promising qubit candidates.

2.6.2 Qubit material host requirements

Therefore, we turn to the search of other QT compatible hosts that offers similar capabilities, but that are more user-friendly. In particular, we need to search for new promising materials that can host a potential point defect. Weber *et al.* [6] proposed in 2010 four criteria that should be met for a solid-state semiconductor material hosting a qubit defect, whereas some of the criteria has already been discussed. An ideal crystalline host should have [6]

- (H1) A wide-band gap to accomodate a deep center.
- (H2) Small spin-orbit coupling in order to avoid unwanted spin flips in the defect bound states.
- (H3) Availability as high-quality, bulk, or thin-film single crystals.
- (H4) Constituent elements with naturally occuring isotopes of zero nuclear spin.

Table (2.1) lists several material host candidates that exhibit promising band gap capable of accommodating a deep level defect. The spin-orbit splitting is an indication of the strength of the spin-orbit interaction, and is taken at the Γ point from the valence-band splitting. A smaller value may indicate less susceptibility to decoherence.

Criterion (H3) is important for scalability and further potential for a large-scale fabrication. The given candidate hosts provided in table (2.1) can all be grown as single crystals, but with varying quality and size.

Material	Band gap E_g (eV)	Spin-orbit splitting Δ_{so} (meV)	Stable spinless nuclear isotopes?
${}^3\text{C-SiC}$	2.39	10	Yes
${}^4\text{H-SiC}$	3.26 [37]	6.8	Yes
${}^6\text{H-SiC}$	3.02	7.1	Yes
AlN	6.13	19 [38]	No
GaN	3.44	17.0	No
AlP	2.45	50 [39]	No
GaP	2.27	80	No
AlAs	2.15	275	No
ZnO	3.44 [40]	-3.5	Yes
ZnS	3.72 [41]	64	Yes
ZnSe	2.82	420	Yes
ZnTe	2.25	970	Yes
CdS	2.48	67	Yes
C (Diamond)	5.5	6	Yes
Si	1.12	44	Yes

Table 2.1: Table taken from Gordon *et al.* [33] that lists a number of tetrahedrally coordinated hosts whose band gaps are larger than 2.0 (eV), and compares it to diamond and Si. All experimental values are from Ref. [30], except for where explicitly cited otherwise.

Normally, nuclear spin is a major source of decoherence for all semiconductor-based quantum technologies. This would exclude the use of all elements in odd groups in the periodic table, since these elements exhibit nonzero nuclear spin. As a result, the spin-coherence time of a paramagnetic deep center [6] might increase. However, nuclear spin can also induce additional quantum degrees of freedom for applications in the right configuration [42]. Therefore, criterion (H4) is not a strict requirement but is a general recommendation for reducing decoherence time.

Weber *et al.* [6] use criteria (H1) – (H4) to specifically find analogies to the NV^{-1} center in other material systems, thus leaving the discussion of other criteria out, such as the choice of crystal system. The atomic configuration and crystal structure of a material strongly influences the properties of a defect, since a defect's orbital and spin structure is dependent on its spatial symmetry [42]. In particular, it is the point group that decides which multiplicity a given energy level should have [43]. A higher defect symmetry group generally facilitates degenerate states, which may give rise to high spin states according to Hund's rules [42, 44]. Inversion symmetry in the host crystal can also be beneficial, resulting in reduced inhomogeneous broadening and

spectral diffusion of optical transitions as a consequence of being generally insensitive to external electric fields [42].

2.6.3 Silicon carbide

Silicon carbide (SiC) is an emerging quantum platform that exists in a wide variety of polytypes, with 3C, 4H and 6H being the most prominent configurations. Several of the polytypes have been demonstrated to host SPEs with a slightly different emitter characteristic, which provides the opportunity to select the desired properties based on the variety of lattice configurations and point defects available [6, 46, 47]. While 3C has a cubic structure, we find 4H in a hexagonal structure with both hexagonal (h) and pseudo-cubic (k) lattice sites. 6H is found in a hexagonal structure with the three orientations that are labelled h, k_1 and k_2 . Importantly, SiC in the three varieties experience wide-band gaps, low spin orbit coupling and stable spinless nuclear isotopes [6, 30, 37], as seen from table 2.1. Furthermore, SiC benefits from mature fabrication on the wafer-scale, which checks the last of the four (H1-H4) QT host requirements, marking it as a suitable quantum material platform.

The most studied emitters in SiC include the carbon antisite-vacancy pair $C_{Si}V_C$ that emits in the red, the silicon vacancy V_{Si} that emits in the near infra-red, and the divacancy ($V_{Si}V_C$) and the nitrogen-vacancy center (N_CV_{Si}) that both emit at near-telecom wavelengths. Thus, the two latter emitters could potentially ease the integration with optic fiber technologies as compared to e.g. the NV^- . Additionally, the four different point defects have all been identified as room-temperature SPEs with demonstrated coherent spin control [48]. Illustrations of several configurations of emitters in 4H-SiC are included in figure 2.7.

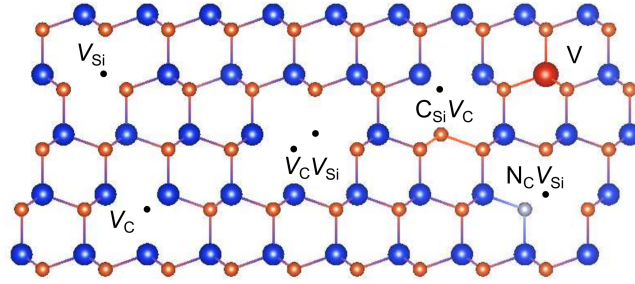


Figure 2.7: Schematic illustration of various point defects in 4H-SiC, where Si atoms are blue while C atoms are orange. The illustration includes the point defects Si vacancy (V_{Si}), C vacancy (V_C), divacancy ($V_{Si}V_C$), carbon antisite-vacancy pair ($C_{Si}V_C$), nitrogen-vacancy (N_CV_{Si}) and the vanadium impurity (V). Figure taken from Ref. [45].

2.6.4 Alternative promising material hosts

Single photon emitters have been observed in other semiconductor materials, however most of the emitters are yet to be identified or are in an early stage of identification. Therefore, specific details about spin- or emission-related structure are yet to be implemented. In this section we will briefly mention recent promising materials for QT.

One immediate potential candidate is silicon, considering the favorable device fabrication processes that are available. It has demonstrated that phosphorous impurities at Si sites can store a quantum state for over 30 seconds, enabling their use in a potential Kane quantum computer [49]. Unfortunately, the P impurity lack any single photon source capabilities. Recently, however, the G-center arising from the carbon-interstitial carbon-substitutional (C_sC_i) complex was identified as an promising SPE candidate with single photon emissions at telecom wavelength [50].

Other materials that emits individual photons have been detected in other wide-band gap semiconductors, including ZnO, ZnS, GaAs, GaN and AlN [49, 51]. Unfortunately, challenges due to the specific materials complicate the implementation of defects for QT. ZnO and ZnS experience a broad emission due to a large photon involvement. GaAs is promising since it has been demonstrated as a SPS, but demonstration of spin manipulation is still in an early phase [51]. GaN and AlN, on the other hand, are more susceptible to a more narrow emission, where room-temperature SPE has been demonstrated for both GaN [52] and wurtzite AlN films [53]. The defect levels for AlN films have been tentatively assigned to the nitrogen-vacancy and divacancy complexes, but they tend to occur too close to the band edges for any SPE [49, 54].

Recent advances in material growth have enabled the use of hole spin-based semiconductors, such as SiGe quantum wells due to their low disorder and large intrinsic spin-orbit coupling strength [55]. Promising materials can also emerge from placing an impurity next to a vacancy. Cation vacancies in possible structures tend to be negatively charged, thus the impurities should act as donors. Therefore, the self-activation center in ZnSe can be a promising defect [6], but is still in the early stage of development.

Two-dimensional materials such as hexagonal boron nitride (h-BN), MoS_2 , WSe_2 and WS_2 are also of interest as quantum platforms [56, 57]. The structure of h-BN exists in single- or multilayers, and it has been demonstrated a broad range of stable room-temperature single-photon emitters [58, 59]. In WSe_2 , $MoSe_2$ and WS_2 , there has been experimentally discovered optical excitation of defects, while also electrical excitation of defects for WS_2 [57]. However, secure identification for the source of the emission is yet to be established [57, 60, 61].

2.6.5 Associated challenges with material host discovery

The idea of finding new potential host candidates to utilise point defects in QT is of a challenging sort. Recall, we have made four criteria that deals with the required (H1) band gaps, (H2) spin-orbit coupling, (H3) availability and (H4) spin-zero isotopes, but we have no knowledge of if there should be more criteria or to what extent a criterion needs to be fulfilled. What we do know is that there are major advantages if materials exhibit properties such as isolation in the lattice and weak electron-photon interaction, however, the process to provide any quantity of measurements are through approximations and material-specific properties. These approximations does not neccessarily capture quantum properties well.

Furthermore, the identified candidates constitutes an immensely selective group of only a handful potential hosts. As an example, most known potential hosts are unary or binary compounds. This is probably due to the increasing complexity dealing with an additional level of interactions in the lattice. Therefore, there are reasons to believe that many potential hosts are yet to be discovered, which serves as a motivation for studies involving exploratory research for new candidates.

Chapter 3

Understanding information from *ab-initio* calculations

In this work we will be looking into information gained by *ab-initio* calculations, which means "from first principles". Initially, the method of the calculations can be regarded as a black box where one provides the structure of a material as input, which the black box in return feed us the outcome in terms of interesting constants or a different structure. In this chapter we provide the necessary knowledge to understand what is happening inside the black box, and the quality of the output. The significance of this will be crystal clear when we later have to interpret information regarding several thousands solid-state systems, and not a single material. In particular, the quality of the output is dependent on being consistent.

The black box is based on a successful theory called density functional theory (DFT), which is an approach for predicting physical properties of solid-state systems. We will only summarize the necessary theory behind density functional theory, leaving most of the quantum-mechanical world untouched. We will start with an explanation of why it is difficult to calculate interactions between particles, and then review key approximations and methods regarding the theory. However, even if density functional theory solves some problems, it also introduce new challenges, which will be thoroughly discussed.

3.1 Introduction to density functional theory

To fully understanding what challenges the density functional theory solves, we will need to investigate how we can calculate the forces acting inside a crystal. Since these forces are happening on a microscopic scale, we will need to utilize the theory of quantum mechanics.

3.1.1 The Schrödinger equation

In principle, we can describe all physical phenomenas of a system with the wavefunction $\Psi(\mathbf{r}, t)$ and the Hamiltonian $\hat{H}(\mathbf{r}, t)$, where \mathbf{r} is the spatial position and t is the time. Unfortunately, analytical solutions for the the time-dependent Schrödinger equation,

$$i\hbar \frac{\partial}{\partial t} \Psi(\mathbf{r}, t) = \hat{H}(\mathbf{r}, t) \Psi(\mathbf{r}, t), \quad (3.1)$$

are extremely rare. More conveniently, we can generate a general wavefunction by a summation of eigenfunctions,

$$\Psi(\mathbf{r}, t) = \sum_{\kappa} c_{\kappa} \psi_{\kappa}(\mathbf{r}, t), \quad (3.2)$$

where c_{κ} is a constant and ψ_{κ} is the κ -th eigenfunction. A general wavefunction does not neccessarily describe stationary states, and consequently does not have distinct energies but is rather represented statistically from the expectation value

$$E = \sum_{\kappa} |c_{\kappa}|^2 E_{\kappa}. \quad (3.3)$$

Solving the Schrödinger equation for a general wavefunction is rather troublesome, but luckily we can use the eigenfunctions instead, transforming equation 3.1 into the time-independent Schrödinger equation for eigenfunctions

$$\hat{H} \psi_{\kappa}(\mathbf{r}) = E_{\kappa} \psi_{\kappa}(\mathbf{r}), \quad (3.4)$$

where E_{κ} is the eigenvalue of the κ -th eigenstate $\psi_{\kappa}(\mathbf{r})$. The eigenfunctions have distinct energies, and the state with the lowest energy is called the ground state. They have the attribute that they are orthogonal and normalized with respect to

$$\langle \psi_{\kappa}(\mathbf{r}) | \psi_{\kappa'}(\mathbf{r}) \rangle = \delta_{\kappa\kappa'}. \quad (3.5)$$

The symmetry of an eigenfunction depends on the symmetry of the potential $V_{\text{ext}}(\mathbf{r})$ and the boundary conditions [62].

3.1.2 The many-particle Schrödinger equation

As we extend the theory to include many-particle systems, we will gradually explain and add the different contributions that make up the many-body Hamiltonian. During this process, we will neglect any external potential applied to the system.

If we place a simple electron with mass m_e in its own system, it will be in possession of kinetic energy. Instead of just one electron, we can place N_e electrons, and they will together have the total kinetic energy

$$T_e = - \sum_{j=1}^{N_e} \frac{\hbar^2 \nabla_j^2}{2m_e}. \quad (3.6)$$

All the electrons are negatively charged, causing repulsive Coulomb interactions between each and every electron, totalling to

$$U_{ee} = \sum_{j=1}^{N_e} \sum_{j' < j} \frac{q^2}{|r_j - r_{j'}|}. \quad (3.7)$$

The summation voids counting each interaction more than once. Simultaneously, we can place N_n nuclei with mass m_n in the same system, accumulating the kinetic energy

$$T_n = - \sum_{a=1}^{N_n} \frac{\hbar^2 \nabla_a^2}{2m_n}. \quad (3.8)$$

As in the example with electrons, the nuclei are also experiencing repulsive interactions between every single nucleus, adding up the total interactions as

$$U_{nn} = \sum_{a=1}^{N_n} \sum_{a' < a} \frac{q^2 Z_a Z_{a'}}{|R_a - R_{a'}|}. \quad (3.9)$$

where Z_a is the atom number of nuclei number a . The system now contains N_e electrons and N_n nuclei, thus we need to include the attractive interactions between the them,

$$U_{en} = - \sum_{j=1}^{N_e} \sum_{a=1}^{N_n} \frac{q^2 Z_a}{|r_j - R_a|}. \quad (3.10)$$

Together, these equations comprise the time-independent many-particle Hamiltonian

$$\begin{aligned} \hat{H} = & - \sum_{j=1}^{N_e} \frac{\hbar^2 \nabla_j^2}{2m_e} - \sum_{a=1}^{N_n} \frac{\hbar^2 \nabla_a^2}{2m_n} + \sum_{j=1}^{N_e} \sum_{j' < j} \frac{q^2}{|r_j - r_{j'}|} \\ & + \sum_{a=1}^{N_n} \sum_{a' < a} \frac{q^2 Z_a Z_{a'}}{|R_a - R_{a'}|} - \sum_{j=1}^{N_e} \sum_{a=1}^{N_n} \frac{q^2 Z_a}{|r_j - R_a|}. \end{aligned} \quad (3.11)$$

A few problems arise when trying to solve the many-particle Schrödinger equation. Firstly, the amount of atoms in a crystal is very, very massive. As an example, we can numerically try to calculate the equation 3.7 for a 1mm^3 silicon-crystal that contains $7 \cdot 10^{20}$ electrons. For this particular problem, we will pretend to use the current fastest supercomputer Fugaku [63] that can calculate 514 TFlops, and we will assume that we need 2000 Flops to calculate each term inside the sum [62], and we need to calculate it $N_e \cdot N_e/2$ times for the (tiny) crystal. The entire electron-electron interaction calculation would take $2.46 \cdot 10^{19}$ years to finish for a tiny crystal. Thus, the large amount of particles translates into a challenging numerical problem.

Secondly, the many-particle Hamiltonian contains operators that has to be applied to single-particle wavefunctions, and we have no prior knowledge of how Ψ depends on the single-particle wavefunctions ψ_κ .

3.1.3 The Born-Oppenheimer approximation

The many-particle eigenfunction describes the wavefunction of all the electrons and nuclei and we denote it as Ψ_κ^{en} for electrons (e) and nuclei (n), respectively. The Born-oppenheimer approximation assumes that nuclei, of substantially larger mass than electrons, can be treated as fixed point charges. According to this assumption, we can separate the eigenfunction into an electronic part and a nuclear part,

$$\Psi_\kappa^{\text{en}}(\mathbf{r}, \mathbf{R}) \approx \Psi_\kappa(\mathbf{r}, \mathbf{R})\Theta_\kappa(\mathbf{R}), \quad (3.12)$$

where the electronic part is dependent on the nuclei. This is in accordance with the assumption above, since electrons can respond instantaneously to a new position of the much slower nucleus, but this is not true for the opposite scenario. From here, one can obtain the electronic and nuclear eigenfunction, with the derivation shown in A.1.

3.1.4 The Hartree and Hartree-Fock approximation

The next question in line is to find a wavefunction $\Psi(\mathbf{r}, \mathbf{R})$ that depends on all of the electrons in the system. The Hartre [62] approximation to this is to assume that electrons can be described independently, suggesting the *ansatz* for a two-electron wavefunction

$$\Psi_\kappa(\mathbf{r}_1, \mathbf{r}_2) = A \cdot \psi_1(\mathbf{r}_1)\psi_2(\mathbf{r}_2), \quad (3.13)$$

where A is a normalization constant. This approximation simplifies the many-particle Schrödinger equation a lot, but comes with the downside that the particles are distinguishable and do not obey the Pauli exclusion principle for fermions.

The Hartree-fock approach, however, overcame this challenge and presented an anti-symmetric wavefunction that made the electrons indistinguishable [1]:

$$\Psi_{\kappa}(\mathbf{r}_1, \mathbf{r}_2) = \frac{1}{\sqrt{2}} \left(\psi_1(\mathbf{r}_1)\psi_2(\mathbf{r}_2) - \psi_1(\mathbf{r}_2)\psi_2(\mathbf{r}_1) \right). \quad (3.14)$$

For systems containing more than one particles, the factor $1/\sqrt{2}$ becomes the Slater determinant and is used to normalize the wave function.

3.1.5 The variational principle

So far, we have tried to make the time-independent Schrödinger equation easier with the use of an *ansatz*, but we do not necessarily have an adequate guess for the eigenfunctions and the ansatz can only give a rough estimate in most scenarios. Another approach, namely the *variational principle*, states that the energy of any trial wavefunction is always an upper bound to the exact ground state energy by definition E_0 .

$$E_0 = \langle \psi_0 | H | \psi_0 \rangle \leq \langle \psi | H | \psi \rangle = E \quad (3.15)$$

This enables a minimization of energy in terms of wavefunction parameters. A more thorough walk-through of the variational principle is included in appendix A.2.

3.2 The density functional theory

Hitherto we have tried to solve the Schrödinger equation to get a ground state wave function, and from there we can obtain ground state properties, such as the ground state total energy. One fundamental problem that exists when trying to solve the many-electron Schrödinger equation is that the wavefunction is a complicated function that depends on $3N_e$ variables¹.

Hohenberg and Kohn [64] showed in 1964 that the ground-state density $n_0(\mathbf{r}) = |\Psi_0(\mathbf{r})|$ determines a general external potential, which includes U_{en} , up to an additive constant, and thus also the Hamiltonian [65]. From another point of view, the theory states that all physical ground-state properties of the many-electron system are unique functionals of the density [62]. A consequence of this is that the number of variables is reduced from $3N_e$ to 3, significantly reducing the computational efforts.

However, the scheme is not without limitations, as the density functional theory (DFT) can only be used to find all the ground-state physical properties

¹not including spin

if the exact functional of the electron density is known. And 57 years after Hohenberg and Kohn published their paper, the exact functional still remains unknown.

We will start this chapter with a brief mention of the Hohenberg-Kohn theorems and its implications, before we delve further into the Kohn-Sham equation. From the

3.2.1 The Hohenberg-Kohn theorems

THEOREM 1. *For any system of interacting particles in an external potential V_{ext} , the density is uniquely determined.*

The theorem can be proved by utilising the variational principle for two different external potentials with the same ground state density. The proof is included in appendix A.3.1.

THEOREM 2. *There exists a variational principle for the energy density functional such that, if n is not the electron density of the ground state, then $E[n_0] < E[n]$.*

From theorem 1, we know that the external potential is uniquely determined by the density, which in turn uniquely determines the ground state wavefunction. Therefore, all other observables of the system are uniquely determined and we can express the energy as function of the density,

$$E[n] = \overbrace{T[n] + U_{ee}[n]}^{F[n]} + U_{en}[n]. \quad (3.16)$$

where $F[n]$ is an universal functional known as the Hohenberg-Kohn functional. The proof for theorem 2 is found in appendix A.3.2.

3.2.2 The Kohn-Sham equation

So far, we have tried to make the challenging Schrödinger equation less challenging by simplifying it, with the last attempt containing the Hohenberg-Kohn's theorems where the theory states that the total ground-state energy can, in principle, be determined exactly once we have found the ground-state density.

In 1965, Kohn and Sham [66] reformulated the Hohenberg-Kohn theorems by generating the exact ground-state density $n_0(\mathbf{r})$ using a Hartree-like total wavefunction

$$\Psi(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_{N_e}) = \psi_1^{\text{KS}}(\mathbf{r}_1) \psi_2^{\text{KS}}(\mathbf{r}_2) \dots \psi_{N_e}^{\text{KS}}(\mathbf{r}_{N_e}), \quad (3.17)$$

where $\psi_j^{\text{KS}}(\mathbf{r}_j)$ are some auxiliary independent single-particle wavefunctions. However, the Kohn-Sham wavefunctions cannot be the correct single-particle wavefunctions since our ansatz implies an exact density

$$n(\mathbf{r}) = \sum_{j=1}^{N_e} |\psi_j^{\text{KS}}(\mathbf{r})|^2. \quad (3.18)$$

Recalling that equation 3.16 describes the total energy as a functional of the density,

$$E[n] = T[n] + U_{ee}[n] + U_{en}[n], \quad (3.19)$$

we try to modify it to include the kinetic energy $T_s[n]$ and the interaction energy $U_s[n]$ of the auxiliary wavefunction, and the denotation s for single-particle wavefunctions.

$$\begin{aligned} E[n] &= T[n] + U_{ee}[n] + U_{en}[n] + (T_s[n] - T_s[n]) + (U_s[n] - U_s[n]) \\ &= T_s[n] + U_s[n] + U_{en}[n] + \underbrace{(T[n] - T_s[n]) + (U_{ee}[n] - U_s[n])}_{E_{xc}[n]} \end{aligned}$$

Here we have our first encounter with the *exchange-correlation energy*

$$E_{xc}[n] = \Delta T + \Delta U = (T[n] - T_s[n]) + (U_{ee}[n] - U_s[n]), \quad (3.20)$$

which contains the complex many-electron interaction. For non-interacting system, $E_{xc}[n]$ is conveniently zero, but in interacting systems it most likely is a complex expression. However, one can consider it as our mission to find good approximations to this term, as the better approximations, the closer we get to the exact expression.

The exact total energy functional can now be expressed as

$$\begin{aligned} E[n] &= \underbrace{\sum_j \int \psi_j^{\text{KS}*} \frac{-\hbar^2 \nabla^2}{2m} \psi_j^{\text{KS}} d\mathbf{r}}_{T_s[n]} + \underbrace{\frac{1}{2} \int \int q^2 \frac{n(\mathbf{r})n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d\mathbf{r}d\mathbf{r}'}_{U_s[n]} \\ &\quad + \underbrace{\int V_{en}(\mathbf{r})n(\mathbf{r})d\mathbf{r}}_{U_{en}[n]} + \underbrace{(T[n] - T_s[n]) + (U_{ee}[n] - U_s[n])}_{E_{xc}[n]}. \end{aligned} \quad (3.21)$$

given that the exchange-correlation functional is described correctly. By utilizing the variational principle, we can now formulate a set of Kohn-Sham single-electron equations,

$$\left\{ -\frac{\hbar^2}{2m_e} \nabla_s^2 + V_H(\mathbf{r}) + V_{j\alpha}(\mathbf{r}) + V_{xc}(\mathbf{r}) \right\} \psi_s^{\text{KS}}(\mathbf{r}) = \epsilon_s^{\text{KS}} \psi_s^{\text{KS}}(\mathbf{r}) \quad (3.22)$$

where $V_{xc}(\mathbf{r}) = \partial E_{xc}[n]/\partial n(\mathbf{r})$ and $V_H(\mathbf{r}) = \int q^2 \frac{n(\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|} d\mathbf{r}'$ is the Hartree potential describing the electron-electron interaction. It is worth to notice that $V_H(\mathbf{r})$ allows an electron to interact with itself, resulting in a self-interaction contribution, however this will be taken care of in V_{xc} .

Finally, we can define the total energy of the system according to Kohn-Sham theory as

$$E[n] = \sum_j \epsilon_j^{KS} - \frac{1}{2} \iint q^2 \frac{n(\mathbf{r})n(\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|} d\mathbf{r}d\mathbf{r}' + E_{xc}[n] - \int V_{xc}(\mathbf{r})n(\mathbf{r})d\mathbf{r}. \quad (3.23)$$

If V_{xc} is exact, and $E[n]$ gives the true total energy, we still do not know if the energy eigenvalues ϵ_s^{KS} are the true single-electron eigenvalues. However, there exists one exception, which is that the highest occupied eigenvalue of a finite system has to be exact if the density is exact.

The only task that is left for us now is to find the exact expression for $E_{xc}[n]$ as a functional of the density $n(\mathbf{r})$. With that expression, we would be able to calculate the total energies of any material, and most likely solve a few of the biggest puzzles in the history of humankind. Unfortunately, the exchange-correlation potential is unknown for most systems.

It is possible to solve the Kohn-Sham equations by applying a self-consistent field method. This is a computational scheme, and for further details one can consult the appendix A.4.

3.2.3 The exchange-correlation energy

There is one scenario for which we can derive the exact expression of the exchange-correlation functional, namely the *homogeneous electron gas* (HEG). However, this has a natural cause, since by definition $n(\mathbf{r})$ is constant for this situation. Given that it is the variations of electron density that are the foundation of material properties, the usefulness of HEG is limited. The *local density approximation* (LDA) is an approximation based on this approach, where the local density is the only variable used to define the exchange-correlation functional. Specifically, we can set the exchange-correlation potential at each position to be the known exchange-correlation potential from homogeneous electron gas at the electron density observed at that position [66]:

$$V_{xc}(\mathbf{r}) = V_{xc}^{\text{electron gas}}[n(\mathbf{r})]. \quad (3.24)$$

This is the simplest and most known approximation to the exchange-correlation functional, and accordingly it has a few drawbacks. One of them is the incomplete cancellation of the self-interaction term, which leads to a repulsion that may cause artificial repulsion between electrons, and hence increased electron delocalization [67]. In addition, LDA has proven challenging to use when

studying atoms and molecules because of their rapidly varying electron densities, however, the LDA is seen as succesful for bulk materials because of the slowly varying electron density [68]. Considering the relatively low computational cost and relatively high accuracy, the LDA overall makes a good model for estimation of the exchange-correlation functional for bulk-materials.

In the light of the merits of the LDA, an extensive search for new approximations was launched. The *generalized gradient approximation* (GGA) is an extension of the LDA, which includes the gradient of the density

$$V_{xc}^{GGA}(\mathbf{r}) = V_{xc} [n(\mathbf{r}), \nabla n(\mathbf{r})] . \quad (3.25)$$

The GGA is a good approximation for the cases where the electron density varies slowly, but faces difficulties in many materials with rapidly varying gradients in the density, causing the GGA to fail. Thus, the annotation *generalized* in GGA is set to include the different approaches to deal with this challenge. Two of the most commonly implemented GGA functionals are the non-empirical approaches Perdew-Wang 91 (PW91) [69] and Perdew-Burke-Ernzerhof (PBE) [70].

Both LDA and GGA are commonly known to severely underestimate the band gaps of semiconductor materials, in addition to incorrectly predicting charge localizations originating from narrow bands or associated with local lattice distortions around defects [71]. The latter limitation is thought to be due to self-interaction in the Hartree potential in equation 3.22.

To improve the accuracy of excited state properties estimations, other methods have been developed. Tran and Blaha (TB) [72] adapted the exchange potential suggested by Becke-Roussel (BR) [73] that leads to band gaps close to experimental values while still using a cheap semilocal method [74]. The modified TB-mBJ version introduces a parameter relative to the BJ version,

$$V_{xc}^{TB-mBJ}(\mathbf{r}) = c \cdot V_{xc}^{BR}(\mathbf{r}) + (3c - 2) \frac{1}{\pi} \sqrt{\frac{5}{6}} \sqrt{\frac{t(\mathbf{r})}{\rho(\mathbf{r})}}, \quad (3.26)$$

where ρ is the electron density, t is the kinetic-energy density, V_{xc}^{BR} is the Becke-Roussel exchange potential and c is the parameter that changes weights of the BJ potential and is fit to experimental data.

Hybrid functionals intermix exact Hartree-Fock exchange with exchange and correlation from functionals based on the LDA or GGA. Hartree-Fock theory completely ignore correlation effects, but account for self-interaction and treats exchange as exact. Since LDA/GGA and Hartree-Fock supplement each other, they can be used as a combination for hybrid-functionals resulting in some cancellation of the self-interaction error. Becke [75] introduced a 50% Hartree-Fock exact exchange and 50% LDA energy functional, while Perdew *et al.* [76] altered it to 25% – 75% and favoring PBE-GGA instead of LDA.

The inclusion of Hartree Fock exchange improves the description of localized states, but requires significantly more computational power for large systems. Another method called the *GW* approximation includes screening of the exchange interaction [77], but has a computational price that does not necessarily defend its use. Thus, the real challenge is to reduce the computational effort while still producing satisfactory results. Heyd *et al.* [78] suggested to separate the non-local Hartree-Fock exchange into a short- and long-range portion, incorporating the exact exchange in the short-range contribution. The separation is controlled by an adjustable parameter ω , which was empirically optimised for molecules to $\omega = 0.15$ and solids to $\omega = 0.11$ and are known as the HSE03 and HSE06 (Heyd-Scuseria-Ernzerhof), respectively [79]. The functionals are expressed as

$$E_{xc}^{HSE} = \alpha E_x^{HF,SR}(\omega) + (1 - \alpha) E_x^{PBE,SR}(\omega) + E_x^{PBE,LR}(\omega) + E_c^{PBE} \quad (3.27)$$

where $\alpha = 1/4$ is the Hartree-Fock mixing constant and the abbreviations SR and LR stands for short range and long range, respectively.

Hence, hybrid-functionals are *semi-empirical* functionals that rely on experimental data for accurate results. They give accurate results for several properties, such as energetics, bandgaps and lattice parameters, and can fine-tune parameters fitted to experimental data for even higher accuracy.

Furthermore, the computational effort required for the hybrid-functionals are significantly larger than for non-empirical functionals such as LDA or GGA. Krukau *et al.* [79] reported a substantial increase in computational cost when reducing the parameter ω from 0.20 to 0.11 for 25 solids, and going lower than 0.11 demanded too much to actually defend its use.

Unfortunately, an area where both GGA and hybrid functionals are reportedly inadequate is in calculating the dispersion interactions [80]. Many implementations have been developed to deal with this vulnerability, with one of them being the non-local van der Waals density functional (vdW-DF) [81]

$$E_{xc} = E_x^{GGA} + E_c^{LDA} + E_c^{nl}, \quad (3.28)$$

where E_x^{GGA} is GGA exchange energy [80].

3.2.4 Limitations of the DFT

If we had known the exact exchange-correlation functional, the density functional theory would yield the exact total energy. Alas, that is not the case and we are bound to use approximations in forms of functionals. What is common for all approximations is that they are specifically designed to solve one given optimization, therefore it is not necessarily one functional that is

consider superior in all fields of interest. One could consider that the hybrid functionals due to a high accuracy overall should be dominant, but that is only if one have the computational capacity required. The accuracy of calculations is dependent on which functional being used, and normally a higher accuracy means the use of a more complex and computationally demanding functional.

Nonetheless, density functional theory is considered a very successful approach and Walter Kohn was awarded the Nobel Price in chemistry in 1998 for his development of the density-functional theory [82]. It is especially regarded as successful in contexts where DFT can make important contributions to scientific questions where it would be essentially impossible to determine through experiments [68]. One can only hope that the future will be as bright as the past, and that this successful theory provides incentives for further growth in the next generation.

Chapter 4

Machine learning

The enormous amount of data generated in the digital world today is beyond comprehension. In 2019, more than 500 hours of video was uploaded to Youtube every second, totalling to over 82 years of content every day¹. In addition, more than 1.5 billion web sites exists².

However, an increasing amount of data comes hand in hand with an increasing demand of knowledge about the data. If we are unable to extract information from the data, the data serves no intention and exists as an excess. Therefore, we need methods to process and automate data analysis, which is what the promises of *machine learning* covers. Machine learning can reveal patterns in data with ease where a human would face difficulties, and use this information to predict or generate new data. Many tools in machine learning is based on probability theory, which can be applied to problems involving uncertainty. Thus, machine learning is also commonly named as *statistical learning* [83].

There are mainly two types of machine learning, either *supervised* or *unsupervised* learning. In unsupervised learning we are given inputs $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$, where \mathbf{x}_i is a training input that has D-dimensions that describes each entry, where each dimension is known as a *feature* or a *descriptor*. The features could be exemplified as height or weight, or it could be something complex that has no practical meaning (at least not to humans). Since there are no features describing what an entry is, it is up to the tools of machine learning to find patterns in the data, and is the essence of unsupervised learning. In the supervised approach, on the other hand, the model tries to learn a mapping from inputs \mathbf{x} to outputs y , given a labeled set of pairs $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$. The set \mathcal{D} is known as the training set, and N is the number of entries. The flexibility of the shape of a feature is also shared with the output. It can in principle

¹Source: <https://www.youtube.com/intl/no/about/press> extracted 15.02.2021

²Source: <https://www.statista.com/chart/19058/how-many-websites-are-there> extracted 29.03.2021

be anything, but it is mostly assumed that the output is either *categorical* or *nominal* restricted by a finite set $y_i \in \{1, \dots, \mathcal{C}\}$. The problem is defined as *classification* if the output is categorical, or *regression* if the output is real-valued [83].

4.1 Supervised learning

Supervised learning applied to classification has as goal to learn the target output $y \in \{1, \dots, \mathcal{C}\}$ from the inputs \mathbf{x} . The number of classes is \mathcal{C} , and depicts if the classification is *binary* ($\mathcal{C} = 2$), *multiclass* ($\mathcal{C} > 2$), or *multi-label* if the class labels are not mutually exclusive (exemplified with the weather can be both sunny and cold at the same time). Normally, classification is used when the problem is formulated as a multiclass classification, and hereon we will adapt to the formulation as well [83].

In order to be able to learn from data, we will need to formulate a function approximate. Assume $y = f(\mathbf{x}) + \epsilon$ for some unknown function f and a random error term ϵ with mean zero. We can then try to approximate f from a labeled training set, which we can use to make the predictions $\hat{y} = \hat{f}(\mathbf{x})$. With the estimated \hat{f} we can make predictions on unlabeled data and achieve a *generalized model*. The estimated function \hat{f} is often considered as a black box, since we are not necessarily interested in the exact shape of the function but rather the predictions.

As simple as the idea behind supervised classification appears, a generalized model remains deeply dependent on the available data. Imagine a training set containing two entries. The first entry is a young and tall person labeled healthy. The other entry is an old and short person labeled sick. The pattern in this simple scenario is abundantly clear, but will face a challenge if it were to predict on a test set containing a person who is young and short. Therefore, it is desirable to compute the probability of an entry belonging to one class. The probability distribution is given by $p(y|\mathbf{x}, \mathcal{D})$, where the probability is conditional on the input vector (test set) \mathbf{x} and the training set \mathcal{D} . If the output is probabilistic, we can compute the estimation to the true label as

$$\hat{y} = \hat{f}(\mathbf{x}) = \operatorname{argmax} f(\mathbf{x})p(y = 1|\mathbf{x}, \mathcal{D}), \quad (4.1)$$

which represents the most probable class label and is known as the *maximum a posteriori* estimate [83].

4.2 Evaluating accuracy of a model

It would be desirable to find one superior model that we could utilize on all types and sizes of datasets. Unfortunately, there is no algorithm that has

this property, since one model might be recognized as best on one particular dataset, while others are far better on other datasets. This is known as the *no free lunch theorem* (Wolpert 1996 [84]). The same goes with evaluating the model - there is no metrics that stand alone as the best metric to evaluate a model. Choosing how to actually evaluate a model can be the most challenging part of a statistical learning procedure.

4.2.1 Bias-variance tradeoff

To illustrate a challenge in choosing the correct parameters, we give an example using the mean squared error (MSE) as a *cost function*, which we want to minimize in order to improve the accuracy of the model [83]. Assume that our data can be represented by

$$\mathbf{y} = f(\mathbf{x}) + \epsilon,$$

where $f(\mathbf{x})$ is an unknown function and ϵ is normally distributed with a mean equal to zero and variance equal to σ^2 . Furthermore, we also assume that the function $f(\mathbf{x})$ can be approximated to a model $\hat{\mathbf{y}}$, where the model is defined by a design matrix \mathbf{X} and parameters β ,

$$\hat{\mathbf{y}} = \mathbf{X}\beta.$$

The parameters β are in turn found by optimizing the mean squared error (MSE) via the cost function

$$C(\mathbf{X}, \beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2 = \mathbb{E} [(\mathbf{y} - \hat{\mathbf{y}})^2].$$

The cost function can be rewritten as

$$\begin{aligned} \mathbb{E} [(\mathbf{y} - \hat{\mathbf{y}})^2] &= \frac{1}{n} \sum_i (f_i - \mathbb{E} [\hat{\mathbf{y}}])^2 + \frac{1}{n} \sum_i (\hat{y}_i - \mathbb{E} [\hat{\mathbf{y}}])^2 + \sigma^2 \\ &= \mathbb{E} [(\mathbf{f} - \mathbb{E} [\hat{\mathbf{y}}])^2] + \text{Var}(\hat{\mathbf{y}}) + \sigma_\epsilon^2 \end{aligned}$$

where $\mathbb{E}[\mathbf{y}] = \mathbf{f}$, $\mathbb{E}[\epsilon] = \mathbf{0}$ and $\text{Var}(\mathbf{y}) = \text{Var}(\epsilon) = \sigma_\epsilon^2$.

The first term on the right hand side is the squared bias, the amount by which the average of our estimate differs from the true mean, while the second term represents the variance of the chosen model. The last term is the variance of the error ϵ , also known as the irreducible error. In general, an estimated function \hat{f} will never be a perfect estimate for f since we can not reduce the error introduced by ϵ . Therefore, any model will always be restricted to an upper bound of accuracy due to the irreducible error.

A model with high variance will typically experience larger fluctuations around the true value, while a model with high bias corresponds to a larger error in the average of estimates. This is schematically visualized as function of model complexity in figure 4.1. If the model is not complex enough due to high bias and low variance, the algorithm can end up not learning the relevant relations between features and output. This is known as *underfitting* [83]. On the other hand, a complex model with low bias and high variance might find trends in random noise from the training data instead of the relevant features, resulting in *overfitting* [83]. An ideal model would be one that simultaneously achieves low variance and low bias. Therefore, we have to do a trade-off between how much bias and variance we would like in the model.

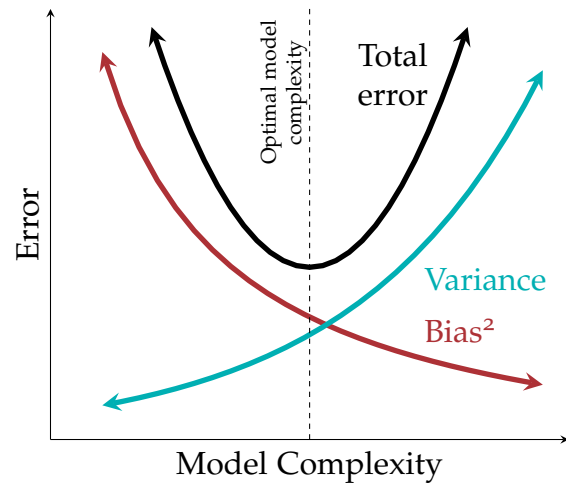


Figure 4.1: A schematic representation of the bias-variance tradeoff as a function of model complexity, adapted from Ref. [85]. The error associated with a model can be decomposed into variance and bias, where a compromise between the lowest bias and variance corresponds to the optimal model complexity.

4.2.2 Accuracy, precision and recall

Given a model that has dealt with the intricacy of increasing complexity, we would like to evaluate the model's output quality. For a binary supervised classification problem we can measure the accuracy by finding how many correct predictions have been made. Prediction accuracy can provide a fine initial analysis, but it has some significant drawbacks seen in unbalanced datasets. This can be easily explained with a dataset consisting of 99 : 1 ratio of class, since just guessing the majority class will result in a very high 99% accuracy. Perhaps it is the 1% that is the most important class, thus the accuracy score severely lacks information for the model.

Therefore, we turn to other evaluation metrics such as a *confusion matrix*. A confusion matrix is a method for measuring the performance of classifiers [83]. It is set up as a table with 4 different categories, where two of the categories are the predicted outcomes of the classifier and the two final categories are the true outcomes. An example of a confusion matrix for a binary classifier is shown in table 4.1.

For the binary confusion matrix there are two possible predicted out-

Table 4.1: A confusion matrix for a binary classifier. The entries true positive and true negative on the diagonal of the matrix are correctly predictions, while false positive and false negative are wrongly made predictions. P and N are the total number positive and negative predictions, respectively. Similarly, P' and N' are the number of true -positive or negative labels, respectively.

		Predicted label		total
		1	0	
Actual label	1	True Positive	False Negative	P'
	0	False Positive	True Negative	N'
total		P	N	

comes, either positive or negative. This gives rise to some terminology.

- True Positive (TP): The classifier correctly predicts a positive event.
- True Negative (TN): The classifier correctly predicts a negative event.
- False Positive (FP): The classifier incorrectly predicts a positive event when the true event was negative.
- False Negative (FN): The classifier incorrectly predicts a negative event when the true event was positive.

From the confusion matrix one can then start estimating the performance of the model, by calculating different factors, such as [83]

- **Sensitivity**, also known as the true negative rate, is the ratio of the number of correct negative examples to the number classified as negative. It is defined as

$$\text{Sensitivity} = \frac{\text{TN}}{\text{TN} + \text{FP}}. \quad (4.2)$$

- **Recall**, also known as the true positive rate, is the ratio of the number of correct positive examples to the number classified as positive. A high recall relates to a low false negative rate, and is defined as

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}. \quad (4.3)$$

- **Precision** is the ratio of correct positive examples to the number of actual positive examples. A high precision relates to a low false positive rate, and is defined as

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}. \quad (4.4)$$

Similar to the bias-variance tradeoff, it is common to compare the recall with the precision to identify the tradeoff for different thresholds. High scores for both reveals that a classifier returns accurate results combined with returning a majority of all positive results.

Sometimes a classifier can have drastically different values for the precision and recall. This leads to another estimator for the performance of a classifier, which is known as the F1-score. The F1-score is defined as the harmonic mean of precision and recall,

$$\text{F1-score} = \frac{2 \cdot \text{Recall} \cdot \text{Precision}}{\text{Recall} + \text{Precision}},$$

and can be used to find a good tradeoff between recall and precision. The highest value of F1-score is 1 and is considered an ideal classifier, while the lowest is 0.

However, the f1 score is insensitive to the number of negative predictions. Therefore, an adjustment of the normal accuracy is in place. The name of this metric is called the balanced accuracy, which equally weights how many true positive and true negative,

$$\text{Balanced accuracy} = \frac{1}{2} (\text{Recall} + \text{Sensitivity}),$$

which makes it particular handy for imbalanced datasets.

We have now only scratched the surface of potential evaluation metrics, and as a final note we would like to emphasize that it is up to the implementer which evaluation metric one should use.

4.2.3 Cross validation

When evaluating different parameters for models, commonly done in a grid-search scheme, there is an abundant risk of performing an overfit to the test set since we can tweak the parameters to a model so it can perform optimally. To solve this problem, we can exclude a part of the dataset as a validation set (in addition to a test set). Therefore, we can train a model on the training set, and evaluate the parameters on the validation set. After a lot of trial and error and the experiment seems successful, we can do one final evaluation on the test set.

Unfortunately, this reduces the number of samples that can be used for training drastically. A fix for this is to apply *cross validation* (CV) [83]. Cross validation is a technique used to evaluate predictive models by partitioning the original sample into a training set to train the model, and a test set to evaluate it.

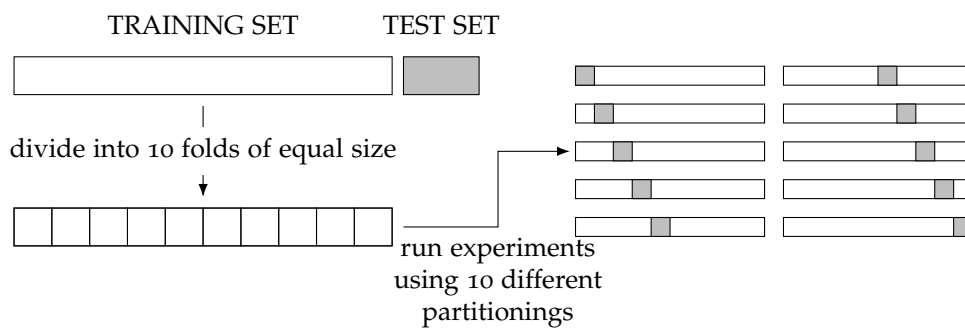


Figure 4.2: A schematic representation of a 10-fold cross-validation scheme.

It is common to apply cross validation into fold, yielding the name of k-fold cross validation. In k-fold cross validation, the training set is partitioned into k equal sized subsamples, as visualized in figure 4.2. Of the k samples, a single sample is used as validation set while the remaining k-1 samples are used as training data. The process is then repeated k-times, such that each of the k subsamples are used as a validation set exactly once. Therefore, all observations are used for both training and validation, and each observation is used for validation exactly once. The k results from the folds can then be averaged to produce an estimate. The subsamples are allowed to have an imbalanced dataset, so that each class is not necessarily represented equally in each fold. Since supervised algorithms tend to weigh each instance equally, this may result in overrepresented classes being favored during training of the model. Even worse could be the result of a fold where one class is not represented at all, resulting in a model that do not learn how to predict a class at all.

To deal with the vulnerability of imbalanced datasets in CV, one can employ a stratified k-fold cross validation technique. Stratification is a process that seeks to ensure that each fold is representative of all classes (also named *strata* in this context) in the data, making each fold having approximately equal class-representation.

4.3 Logistic regression

Logistic regression, or *logit*, is considered a *soft* classification algorithm, which means that an output of the algorithm is considered to be categorical instead of numerical.

Assume we have a dataset with $\mathbf{x}_i = x_{i1}, x_{i2}, \dots, x_{ip}$ input data, where we have p predictors for each corresponding output data y_i . The outcomes y_i are discrete and can only take certain values or classes. In our case we have two classes with y_i either being equal to 0 or 1. Therefore, the probability that a datapoint belongs to either class can be given by the Sigmoid function,

$$p(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}.$$

Furthermore, we have the parameters $\beta = \beta_1, \beta_2, \dots, \beta_p$ of our fitting of the Sigmoid function, where the probabilities are defined as

$$p(y_i | \mathbf{x}_i \beta) = \frac{e(\mathbf{x}_i \beta)}{1 + e(\mathbf{x}_i \beta)}.$$

The goal of logistic regression is then to correctly predict the category of a given dataset, which has different outcomes, by using an optimal parameter β that maximises the probability of seeing the observed data. How we find the parameters $\beta = \beta_1, \beta_2, \dots, \beta_p$ of the model, is to use the principle of *maximum likelihood estimation* (MLE),

$$P(\beta) = \prod_{i=1}^n [p(y_i = 1 | \mathbf{x}_i \beta)]^{y_i} [1 - p(y_i = 1 | \mathbf{x}_i \beta)]^{1-y_i},$$

where we obtain the log-likelihood function, which is easier to work with, since the log-likelihood turns the exponentials into summations,

$$C(\beta) = \sum_{i=1}^n \left(y_i (\mathbf{x}_i \beta) - \log (1 + \exp \{ (\mathbf{x}_i \beta) \}) \right).$$

Finally, choose our cost function for the *cross entropy*, which is defined as the negative log-likelihood,

$$C(\beta) = - \sum_{i=1}^n \left(y_i(\mathbf{x}_i\beta) - \log(1 + \exp\{\mathbf{x}_i\beta\}) \right)$$

To maximize the accuracy and precision of the logistic regression model, we need to find the optimal parameters β by minimising the cross entropy.

4.3.1 Stochastic gradient descent

One common numerical method for finding the minimum of a function is *stochastic gradient descent* (SGD). The fundamental idea of SGD comes from the observation that the cost function can be written as a sum over n data points $\{\mathbf{x}_i\}_{i=1}^n$,

$$C(\beta) = \sum_{i=1}^n c_i(\mathbf{x}_i, \beta).$$

We can compute the gradient as

$$\nabla_{\beta} C(\beta) = \sum_{i=1}^n \nabla_{\beta} c_i(\mathbf{x}_i, \beta)$$

Then, it is possible to introduce a randomness by only taking the gradient on a small interval of the data, called a minibatch. With n total data points, and M datapoints per minibatch, the number of minibatches is then $\frac{n}{M}$.

The idea is now to approximate the gradient by replacing the sum over all data points with a sum over the data points in one of the minibatches picked at random in each gradient descent step,

$$\nabla_{\beta} C(\beta) = \sum_{i=1}^n \nabla_{\beta} c_i(\mathbf{x}_i, \beta) \rightarrow \sum_{i \in B_k} \nabla_{\beta} c_i(\mathbf{x}_i, \beta),$$

where B_k is the set of all minibatches, with $k = 1, \dots, \frac{n}{M}$. One step of gradient descent is then defined by

$$\beta_{j+1} = \beta_j - \gamma_j \sum_{i \in B_k} \nabla_{\beta} c_i(\mathbf{x}_i, \beta)$$

where k is picked at random with equal probability from $[1, \frac{n}{M}]$ and γ_j is the step length. An iteration over the number of minibatches ($\frac{n}{M}$) is commonly referred to as an epoch. Thus, it is typical to choose a number of epochs and for each epoch iterate over the number of minibatches.

4.4 Decision trees

Classification and regression trees (CART), also called decision trees, are one of the more basic supervised algorithms. They can be used for both regression and classification tasks, but we will for the relevancy of this work provide a special emphasis to classification trees.

The idea behind decision trees is to find the features that contain the most information regarding the target, and then split up the dataset along the values of these features. This feature selection enables the target values for the resulting underlying dataset to be as *pure* as possible, which means the dataset only contains one class [83]. The features that can reproduce the best target features are normally said to be the most informative features.

A decision tree can be divided into a *root node*, *interior nodes*, and the final *leaf nodes*, commonly known as *terminal nodes*. The nodes are connected by *branches*. The decision tree is able to learn an underlying structure of the training data and can, given some assumptions, make predictions on unseen observations. These predictions are based on the information stored in the leaf nodes in the tree.

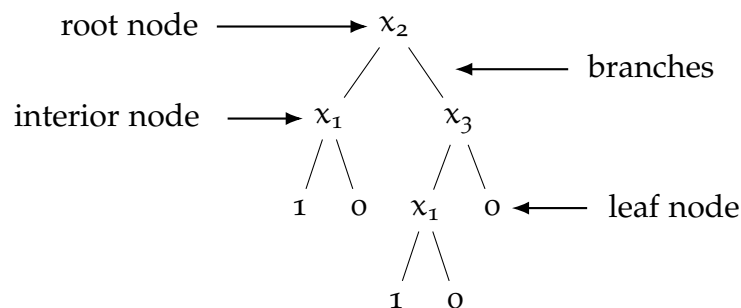


Figure 4.3: A schematic representation of a binary classification tree, which consists of three nodes that contain information of the features x_1 , x_2 and x_3 .

The process behind a decision tree can be seen as a top-down approach. First, we make a leaf provide the classification of a given instance. Then, a node specifies a test of some attribute of the instance, while a branch corresponds to a possible value of an attribute. Subsequently, the instance move down the tree branch corresponding to the value of the attribute. Then the steps can be repeated for a new subtree rooted at the new node.

A classification tree differs from a regression tree by the response of the prediction, since it produces a qualitative response rather than a quantitative one. The response is given by the most commonly occurring class of training observations specified by the attribute of the node. A schematic representation of a classification tree is visualised in figure 4.3.

4.4.1 Growing a classification tree

In growing a classification tree, a process called recursive binary splitting is applied. This involves two steps:

1. Split the set of possible values (x_1, x_2, \dots, x_p) into J distinct non-overlapping regions R_1, R_2, \dots, R_J .
2. If an observation falls within the region R_J , we make the prediction given by the most commonly occurring class of training observations in R_J .

The computational aspect of recursively doing this for every possible combination of features does not defend its use, and therefore the common strategy is to use a top-down approach. Binary splitting begins at the top of the tree and consecutively splits the *predictor space*, which is a space that describes all possible combinations of the features in the dataset. This is indicated by two new branches further down the tree. It should be noted that the top-down approach is a greedy approach, since the best split is made at each step of the tree-growing process, instead of trying to pick a split that will lead to a better tree in a future step.

We can define a *probability density function* (PDF) p_{mk} that represents the number of observations k in a region R_m with N_m observations. This likelihood function can be represented in terms of observations of a class in region R_m as

$$p_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i = k), \quad (4.5)$$

where the *indicator* I function equals zero if we misclassify and one if we classify correctly. Therefore, we can define the splitting of the nodes by the misclassification error

$$m_e = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i \neq k) = 1 - p_{mk}. \quad (4.6)$$

However, other methods exists such as the Gini index

$$g = \sum_{k=1}^K p_{mk}(1 - p_{mk}) \quad (4.7)$$

and the information entropy

$$s = - \sum_{k=1}^K p_{mk} \log p_{mk}. \quad (4.8)$$

The two latter approaches are more sensitive to node purity than the misclassification error, i.e. only containing one class, and are in general preferred [83] for splitting of the nodes in a decision tree.

4.4.2 Classification algorithm

The CART algorithm splits the data set in two subsets using a single feature k and a threshold t_k . The pair of quantities (k, t_k) that constitute the purest subset using the Gini factor G results in the cost function

$$C(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}, \quad (4.9)$$

where G_{left} (G_{right}) measures the impurity of the left (right) subset and m_{left} (m_{right}) is the number of instances on the left (right) subset. The algorithm tries to minimize the cost function to find the pair (k, t_k) by splitting the training set in two, and then following the same logic for the next subsets. It will continue to do this recursively until it reaches the maximum depth hyperparameter, or if the next split does not reduce impurity.

4.4.3 Pruning a tree

A decision tree has the ability to turn into a very complex model, making it prone to overfitting. Pre-pruning is a method that stops the growing of a tree if the decrease in error is not sufficient to justify an increasingly complex model by adding an extra subtree. However, this method should not be implemented for models with a large number of features, since features with small predictive powers might be extensively removed which might result in a tree without any splits at all [83]. Post-pruning, or just pruning, is the standard method which involves growing the tree to full size, and then prune the tree by cutting branches. To determine how much to prune it, we can use a cross-validated scheme to evaluate the amount of terminal nodes that has the lowest error.

4.4.4 Pros and cons of decision trees

Decision trees have several clear advantages compared to other algorithms. They are easy to understand and can be visualised effortlessly for small trees. The algorithm is completely invariant to scaling of the data since each feature is processed separately. Additionally, decision trees can handle both continuous and categorical data and can model interactions between different descriptive features.

As auspicious the advantages of decision trees seems, they are inevitably prone to overfitting and hence does not generalize well. Even with pre-pruning, post-pruning and setting a maximum depth of terminal nodes, the algorithm is still prone to overfit [85]. Another important issue concerns training on unbalanced datasets where one class occurs more frequently than other classes, since this will lead to biased trees because the algorithm will

favor the more occurring class. Furthermore, small changes in the data may lead to a completely different tree. Many of these issues can be addressed by using ensemble methods such as either bagging, random forest, or boosting, and can result in a solid improvement of the predictive performance of trees.

4.5 Ensemble methods

By using a single decision tree, we often end up with an overfitted model that possess a high variance. Luckily, we can apply methods that aggregate different machine learning algorithms to reduce variance. If each of the algorithms get slightly different results, as they learn different parts of the data, we can combine the results into something that is better than any one algorithm alone. These approaches fall under the category of ensemble methods, and will be elaborated upon in this section.

4.5.1 Bagging

Bootstrap aggregation, or just *bagging*, is an ensemble method that involves averaging many estimates [83]. If we have M trained trees on different subsamples of the data, chosen randomly, we can compute the ensemble

$$f(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B f_b(\mathbf{x}), \quad (4.10)$$

where f_b is the b 'th tree. Simply re-running the same algorithm on different subsamples can result in a small variance reduction compared to a single tree due to highly correlated predictors, which showcase the need for better approaches.

Random forests provide an improvement of normal bagged trees by choosing a random sample of m predictors as split candidates from the full set of p predictors. The split is restricted in choosing only one of the m predictors, which are normally chosen as either $m \approx \sqrt{p}$ or $m \approx \log p$. This means that at each split in a tree, the algorithm is restricted to a very small portion of the

available predictors.

Algorithm 1: Random forest algorithm.

```

for For  $b = 1 : B$  do
    Draw a bootstrap sample from the training data;
    Select a tree  $T_b$  to grow based on the bootstrap data;
    while node size smaller than maximum node size do
        Select  $m \leq p$  variables at random from  $p$  predictors;
        Pick the best split point among the  $m$  features using CART
        algorithm and create a new node;
        Split the node into daughter nodes;
    end
end
Output the ensemble of trees  $\{T_b\}_{b=1}^B$  and make predictions

```

By introducing randomness into the model, we arrive at a suprisingly capable model that has a high predictive accuracy [86]. This can be exemplified by supposing that there is one strong predictor in a dataset, together with several other fairly strong predictors. Most of the trees will use this strong predictor at the top split, which means that the bagged trees will look quite similar to each other and will have highly correlated predictions.

However, even with higher prediction accuracy, it comes as a compromise since we lose the easy ability of model interpretation. A single tree can be easy to understand, but interpretation of a huge jungle of trees does not necessarily seem appealing for even an experienced data scientist. Furthermore, a random forest does not substantially reduce the variance as averaging many uncorrelated trees would do, as we will soon find out.

4.5.2 Boosting

Boosting is an ensemble method that fits an additive expansion in a set of elementary basis functions [83]. The basic idea is to combine several weak classifiers, that are only just better than a random guess, in order to create a good classifier. This can be done in an iterative approach where we apply a weak classifier to modify the data. For each iteration, we make sure to weigh the observations that are misclassified with a factor. The method is known as adaptive boosting, since the algorithm is able to adapt during the learning process.

In *forward stagewise additive modeling* we want to find an adaptive model

$$f_M(\mathbf{x}) = \sum_{m=1}^M \beta_m G_m(\mathbf{x}; \gamma_m), \quad (4.11)$$

where β_m are expansion parameters that will be determined in a minimization process, and $b(\mathbf{x}; \gamma_m)$ are functions of the multivariable parameter \mathbf{x} that are described by the parameters γ_m . We will in this example consider a binary classification problem with the outcomes $\gamma_i \in \{-1, 1\}$ where $i = 0, 1, 2, \dots, n-1$ are the set of observables. The predictions are produced by the classification function $G(\mathbf{x})$. The error rate of the training sample is given as

$$\overline{\text{err}} = \frac{1}{n} \sum_{i=0}^{n-1} I(\hat{y}_i \neq G(\mathbf{x}_i)). \quad (4.12)$$

After defining a weak classifier, we can apply it iteratively to repeatedly modified versions of the data producing a sequence of different weak classifiers $G_m(\mathbf{x})$. The iterative procedure can be defined as

$$f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \beta_m G_m(\mathbf{x}), \quad (4.13)$$

where the function $f_M(\mathbf{x})$ will be expressed in terms of

$$G(\mathbf{x}) = \text{sign} \sum_{i=1}^M \alpha_m G_m(\mathbf{x}), \quad (4.14)$$

where α_m is the weight that describes the contribution from the weak classifier $G_m(\mathbf{x})$. The main idea is that we do not go back and adjust earlier parameters, which is why this is called *forward* stagewise additive modeling.

We can demonstrate a binary classification example using the exponential cost function that leads to the *discrete AdaBoost* algorithm [87] at step m ,

$$C(\mathbf{y}, \mathbf{f}) = \sum_{i=0}^{n-1} w_i^m \exp(-\hat{y}_i \beta G(\mathbf{x}_i)), \quad (4.15)$$

where $w_i^m = \exp(-\hat{y}_i f_{m-1}(\mathbf{x}_i))$ is the weight of the corresponding observable i . We can optimize G for any $\beta > 0$ with

$$G_m(\mathbf{x}) = \text{sign} \sum_{i=0}^{n-1} w_i^m I(\hat{y}_i \neq G(\mathbf{x}_i)). \quad (4.16)$$

This is the classifier that minimize the weighted error rate in predicting y . Furthermore, we can rewrite the cost function to

$$C = (\exp(\beta) - \exp(-\beta)) \sum_{i=0}^{n-1} w_i^m I(\hat{y}_i \neq G(\mathbf{x}_i)) + \exp(-\beta) \sum_{i=0}^{n-1} w_i^m. \quad (4.17)$$

Substituting G_m into C and solving for β , we obtain

$$\beta_m = \frac{1}{2} \log \frac{1 - \overline{\text{err}}}{\overline{\text{err}}}, \quad (4.18)$$

with the error redefined as

$$\overline{\text{err}} = \frac{1}{n} \frac{\sum_{i=0}^{n-1} w_i^m I(\hat{y}_i \neq G_m(\mathbf{x}_i))}{\sum_{i=0}^{n-1} w_i^m}. \quad (4.19)$$

Finally, this leads to an update of $f_m(\mathbf{x})$ as defined in equation 4.13 and the weights at the next iteration becomes

$$w_i^{m+1} = w_i^m \exp(-\hat{y}_i \beta_m G_m(\mathbf{x}_i)). \quad (4.20)$$

With the above definitions, we can define the discrete Adaboost algorithm in Algorithm 2.

Algorithm 2: Discrete Adaboost algorithm.

Initialize weights $w_i = 1/n$, $i = 0, \dots, n-1$, such that $\sum_{i=0}^{n-1} w_i = 1$;

for $m = 1 : M$ **do**

 Fit the classifier $f_m(\mathbf{x}) \in \{-1, 1\}$ using weights w_i on the training data;

 Compute the error $\overline{\text{err}} = \frac{1}{n} \frac{\sum_{i=0}^{n-1} w_i^m I(\hat{y}_i \neq G_m(\mathbf{x}_i))}{\sum_{i=0}^{n-1} w_i^m}$;

 Define a quantity $\alpha_m = \log [(1 - \overline{\text{err}}_m)/\overline{\text{err}}_m]$;

 Set new weights to $w_i \leftarrow w_i \exp(\alpha_m I(y_i \neq G(\mathbf{x}_i)))$;

end

 Compute the new classifier $G(\mathbf{x}) = \sum_{i=0}^{n-1} \alpha_m I(y_i \neq G(\mathbf{x}_i))$;

It is possible to apply different cost functions resulting in a variety of boosting algorithms. AdaBoost is an example with the cost function in equation 4.17. But instead of deriving new versions of boosting based on different cost functions, we can find one generic method. This approach is known as *gradient boosting* [88]. Initially, we want to minimize

$$\hat{\mathbf{f}} = \underset{\mathbf{f}}{\text{argmin}} L(\mathbf{f}), \quad (4.21)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))$ are the parameters of the models, and L is a chosen loss function.

This can be solved stagewise, using an approach named *gradient descent*. At step m , let \mathbf{g}_m be the gradient evaluated at $f(\mathbf{x}_i) = f_{m-1}(\mathbf{x}_i)$:

$$\mathbf{g}_m(\mathbf{x}_i) = \left[\frac{\partial L(y_i, f(\mathbf{x}_i))}{\partial f(\mathbf{x}_i)} \right]_{f(\mathbf{x}_i)=f_{m-1}(\mathbf{x}_i)}. \quad (4.22)$$

Then we can update

$$\mathbf{f}_m = \mathbf{f}_{m-1} - \rho_m \mathbf{g}_m, \quad (4.23)$$

where ρ_m is the step length and can be found by approximating the real function

$$\mathbf{h}_m(\mathbf{x}) = -\rho \mathbf{g}_m(\mathbf{x}). \quad (4.24)$$

So far, this only optimizes f at a fixed set of points, but we can modify it by fitting a weak classifier to approximate the negative gradient. Additionally, we add a step length parameter $0 < \nu < 1$ to perform partial updates, also known as *shrinking* [83]. The gradient boost algorithm is shown in Algorithm 3.

Algorithm 3: Gradient boost algorithm.

Initialize the estimate $\mathbf{f}_0(\mathbf{x})$;

for $m = 1 : M$ **do**

 Compute the negative gradient vector $\mathbf{u}_m = -\partial C(\mathbf{y}, \mathbf{f}) / \partial \mathbf{f}(\mathbf{x})$ at

$\mathbf{f}(\mathbf{x}) = \mathbf{f}_{m-1}$;

 Fit the base learner to the negative gradient $\mathbf{h}_m(\mathbf{u}_m, \mathbf{x})$;

 Update the estimate $\mathbf{f}_m(\mathbf{x}) = \mathbf{f}_{m-1}(\mathbf{x}) + \nu \mathbf{h}_m(\mathbf{u}_m, \mathbf{x})$;

end

Output the final estimation $\mathbf{f}_M(\mathbf{x}) = \sum_{m=1}^M \nu \mathbf{h}_m(\mathbf{u}_m, \mathbf{x})$

4.6 Dimensionality reduction

Supervised learning introduces models that can be easy to understand, visualize and has well-defined tools and models. However, a dataset can be tedious to work with due to a large amount of descriptors. These descriptors may also be correlated, which means that no new information will be learned from a correlated feature and therefore could be disregarded. Furthermore, a large dataset pose a computational challenge and a reduction in descriptors could potentially reduce the computational time and effort required for any data analysis. Therefore, it would be beneficial to apply a method that finds correlated descriptors and reduce dimensionality of a dataset. This is the idea of *principal component analysis* (PCA).

4.6.1 Principal component analysis

Principal component analysis is an algorithm that tries to find a low-dimension representation of a dataset that contains as much of the variance in the data as possible [83, 89]. Each of the dimensions found by PCA are a linear combination of the features in the dataset, and are known as *principal components*.

We can write the design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$, with p features and n entries, in terms of its column vectors as

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_0 & \mathbf{x}_1 & \mathbf{x}_2 & \dots & \dots & \mathbf{x}_{p-1} \end{bmatrix}, \quad (4.25)$$

with a given vector

$$\mathbf{x}_i^T = \begin{bmatrix} x_{0,i} & x_{1,i} & x_{2,i} & \dots & \dots & x_{n-1,i} \end{bmatrix}. \quad (4.26)$$

Then we can compute the *covariance matrix* of the design matrix \mathbf{X} , which is a measurement of the joint variability of the p features in \mathbf{X} . The covariance is defined as

$$\text{cov}[\mathbf{v}, \mathbf{u}] = \frac{1}{n} \sum_{i=0}^{n-1} (v_i - \bar{v})(u_i - \bar{u}), \quad (4.27)$$

where \mathbf{v} and \mathbf{u} are two vectors with n elements each. The covariance matrix is defined by applying the covariance for every pairwise feature, resulting in a $p \times p$ matrix. On the diagonal, the covariance of two equal features becomes the variance of one,

$$\text{cov}[\mathbf{u}, \mathbf{u}] = \frac{1}{n} \sum_{i=0}^{n-1} (u_i - \bar{u})(u_i - \bar{u}) = \text{var}[\mathbf{u}]. \quad (4.28)$$

The covariance can take values between $\pm\infty$, which give rise to computational issues due to loss of numerical precision. Therefore, we scale the covariance matrix by the variance. This is known as the correlation function

$$\text{corr}[\mathbf{u}, \mathbf{v}] = \frac{\text{cov}[\mathbf{u}, \mathbf{v}]}{\sqrt{\text{var}[\mathbf{u}]\text{var}[\mathbf{v}]}}. \quad (4.29)$$

Since all values are inbetween -1 and 1 , we avoid any loss of numerical precision. The resulting covariance matrix $\mathbf{C} \in \mathbb{R}^{p \times p}$ becomes

$$\mathbf{C}[\mathbf{x}] = \begin{bmatrix} \text{var}[\mathbf{x}_0] & \text{cov}[\mathbf{x}_0, \mathbf{x}_1] & \text{cov}[\mathbf{x}_0, \mathbf{x}_2] & \dots & \dots & \text{cov}[\mathbf{x}_0, \mathbf{x}_{p-1}] \\ \text{cov}[\mathbf{x}_1, \mathbf{x}_0] & \text{var}[\mathbf{x}_1] & \text{cov}[\mathbf{x}_1, \mathbf{x}_2] & \dots & \dots & \text{cov}[\mathbf{x}_1, \mathbf{x}_{p-1}] \\ \text{cov}[\mathbf{x}_2, \mathbf{x}_0] & \text{cov}[\mathbf{x}_2, \mathbf{x}_1] & \text{var}[\mathbf{x}_2] & \dots & \dots & \text{cov}[\mathbf{x}_2, \mathbf{x}_{p-1}] \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \text{cov}[\mathbf{x}_{p-1}, \mathbf{x}_0] & \text{cov}[\mathbf{x}_{p-1}, \mathbf{x}_1] & \text{cov}[\mathbf{x}_{p-1}, \mathbf{x}_2] & \dots & \dots & \text{var}[\mathbf{x}_{p-1}] \end{bmatrix}, \quad (4.30)$$

for all vectors \mathbf{x}_i where $i = 0, 1, \dots, p-1$. The correlation matrix $\mathbf{K}[\mathbf{x}]$ becomes

$$\mathbf{K}[\mathbf{x}] = \begin{bmatrix} 1 & \text{corr}[\mathbf{x}_0, \mathbf{x}_1] & \text{corr}[\mathbf{x}_0, \mathbf{x}_2] & \dots & \dots & \text{corr}[\mathbf{x}_0, \mathbf{x}_{p-1}] \\ \text{corr}[\mathbf{x}_1, \mathbf{x}_0] & 1 & \text{corr}[\mathbf{x}_1, \mathbf{x}_2] & \dots & \dots & \text{corr}[\mathbf{x}_1, \mathbf{x}_{p-1}] \\ \text{corr}[\mathbf{x}_2, \mathbf{x}_0] & \text{corr}[\mathbf{x}_2, \mathbf{x}_1] & 1 & \dots & \dots & \text{corr}[\mathbf{x}_2, \mathbf{x}_{p-1}] \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \text{corr}[\mathbf{x}_{p-1}, \mathbf{x}_0] & \text{corr}[\mathbf{x}_{p-1}, \mathbf{x}_1] & \text{corr}[\mathbf{x}_{p-1}, \mathbf{x}_2] & \dots & \dots & 1 \end{bmatrix}. \quad (4.31)$$

The covariance matrix can be rewritten as a function of the design matrix,

$$\mathbf{C}[\mathbf{x}] = \mathbb{E}[\mathbf{X}\mathbf{X}^T] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}^T], \quad (4.32)$$

where $\mathbb{E}[\mathbf{X}\mathbf{X}^T]$ is the expectation value, and assuming we have normalized the data such that $\mathbb{E}[\mathbf{X}] = 0$, we can remove the last term.

Further on, we assume that we can do apply a number of orthogonal transformations by some orthogonal matrices $\mathbf{S} = [\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_{p-1}] \in \mathbb{R}^{p \times p}$ with the column vectors $\mathbf{s}_i \in \mathbb{R}^p$. Additionally, we assume that there is a transformation

$$\mathbf{C}[\mathbf{y}] = \mathbf{S}\mathbf{C}[\mathbf{x}]\mathbf{S}^T = \mathbb{E}[\mathbf{S}\mathbf{X}\mathbf{X}^T\mathbf{S}^T], \quad (4.33)$$

such that the new matrix $\mathbf{C}[\mathbf{y}]$ is diagonal with elements $[\lambda_0, \lambda_1, \lambda_2, \dots, \lambda_{p-1}]$. By multiplying with \mathbf{S}^T , we arrive at the given eigenvalue number i of the covariance matrix that

$$\mathbf{s}_i^T \lambda_i = \mathbf{C}[\mathbf{x}] \mathbf{s}_i^T. \quad (4.34)$$

Dimensions with large eigenvalue have a large variation and can therefore be used to find features with useful information since we multiply the eigenvalue with the eigenvectors. When the eigenvalues are small, it means that the eigenvectors shrink accordingly and there is a small variation in these specific features [90].

So far, we have been leading up to the classical PCA theorem. Assume that the data is represented as in equation 4.25 with $\mathbb{E}[\mathbf{X}] = 0$, and assume that there exists an orthogonal transformation $\mathbf{W} \in \mathbb{R}^{p \times p}$. We can then define the reconstruction error

$$J(\mathbf{W}, \mathbf{Z}) = \frac{1}{n} \sum_i (\mathbf{x}_i - \bar{\mathbf{x}}_i)^2, \quad (4.35)$$

with $\bar{\mathbf{x}}_i = \mathbf{W}\mathbf{z}_i$, where \mathbf{z}_i is a column vector with dimension \mathbb{R}^n of the matrix $\mathbf{Z} \in \mathbb{R}^{p \times n}$.

The PCA theorem states that minimizing the above reconstruction error corresponds to setting $\mathbf{W} = \mathbf{S}$, which is the orthogonal matrix that diagonalizes the covariance matrix [83]. The optimal number of features that correspond to the encoding is given by the set of vectors \mathbf{z}_i with at most l vectors. This is defined as the orthogonal projection of the data onto the columns spanned by the eigenvectors of the covariance matrix. Instead of using the covariance matrix, it is preferable to use the correlation matrix to avoid loss of numerical precision. Additionally, it is important to mention that the covariance matrix is sensitive to the standardization of variables, which is why one should always remember to center the data around before applying PCA. We recommend the reader to read Ref. [83] p. 387 for proof of the classical PCA theorem, as we will not elaborate any further. The algorithm for PCA is shown in Algorithm 4.

Algorithm 4: Principal component analysis algorithm.

Set up the design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ with p features and n entries;
 Center the data by subtracting the mean value for each column;
 Compute the covariance matrix $\mathbb{E}[\mathbf{X}\mathbf{X}^T]$;
 Find the eigenpairs of \mathbf{C} with eigenvalues $[\lambda_0, \lambda_1, \dots, \lambda_{p-1}]$ and eigenvectors $[\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_{p-1}]$;
 Order the eigenvalues, and therefore also the eigenvectors, in descending order. Keep only those l eigenvalues larger than a selected threshold value.

Instead of choosing an arbitrary number of dimensions to reduce down to, it is common to choose the number of dimensions that accumulate a sufficient amount of variance. However, it remains a subjective analysis in how many principal components one should include as it will depend on both the specific application and specific data set. If it is impossible to give a motivation for reducing a large dataset to just two or three principal components, there might still be a reason for why to apply PCA to a dataset. PCA can be applied as a preprocessing method to reduce the dimensionality of a dataset, and therefore might drastically improve the efficiency of further supervised learning approaches.

4.7 Practical challenges associated with machine learning

So far, we have covered substantially researched topics such as dimensionality reduction, supervised algorithms and metric evaluation. However, there exists parts of machine learning that does not necessarily get as much attention, but yet are crucial for the objective of machine learning. In this section,

we will briefly mention both known and unknown challenges that are part of building a machine learning model.

The initial phase consists of gathering information systematically. This could be perhaps the most time-consuming part of the entire process, motivated by questions such as how much data is necessary. The answer to this question is as vague as the question itself, since there is no lower or upper bound but rather a general recommendation that the more data the better. Additionally, we should have an hypothesis that we are collecting descriptors of something that can explain the objective of the entire machine learning process. Indeed, the promises of machine learning are limited to data containing good descriptors. For a supervised learning algorithm, it is necessary to have one descriptor for the training data that contains information about what should be learned.

Thereafter follows an analysis of the data quality, often called *pre-processing*. This includes identifying outliers and finding out what to do about any potential missing value in the data. Normally, solutions such as removing outliers and filling any missing value with either the mean, median or zero are applied. The data is also required to be transformed into continuous or categorical values. For the latter case we can carry out a one-hot encoding to ensure that any algorithm does not assume one category being more important to another due to a larger number. Furthermore, it could be necessary to scale the data and reduce dimensionality, with the motivation discussed in section 4.6.1.

If the algorithm has not been chosen yet, this is the time to do so. A clever first-hand approach is to apply a simple algorithm that is not computational demanding to see how it performs on a subset of the data. If the performance is satisfactory, any implementation of a more sophisticated algorithm could be redundant.

Next, the search for optimal hyperparameters while maintaining a generalized model can pose a challenge, but is achievable. It is popular to apply a cross validation during this process with different evaluation metrics, as discussed in section 4.2.

Eventually, with a chosen algorithm and its optimal hyperparameters, we train the algorithm on the entire preprocessed training data and then perform predictions on unseen data. To avoid any bias in the predictions, it is crucial that predictions on unseen data is done only once. The reason for this is that we do not want to optimize any model for the actual test data, since this would reduce the generalization and increase the bias.

Part III

Methodology and implementation

Chapter 5

Material Science Databases

There are multiple different databases for material science discovery available for every day use, some of them completely open-source while others are commercial. This chapter will give a brief overview of databases available for computational material science, and will serve as a toolbox for how to request information and what kind of python packages exist to process that information.

5.1 Fundamentals of a database

A quick search online will reveal the tremendous escalation of effort for big-data driven material science the last few years, resulting in several databases that stores ab-initio calculation details and results. We will here distinguish between a *cloud service*, which is a place to store independent databases for research and commercial purposes, and a *database*, which is an organized collection of structured information. As an example, a cloud service can store several databases, but a database cannot host a cloud service.

To limit the quest of databases, we have restricted the search for databases and cloud services to include inorganic compounds obtained experimentally or by first-principles calculations, in particular DFT-calculations using *Vienna ab initio simulation package* (VASP) [91]. VASP is a software for atomic scale materials programming. Table 5.2 and 5.3 shows a selection of databases and cloud services that meets the given criteries, respectively.

5.1.1 API and HTTP requests

To extract information from a database it is convenient to interact through an *API* (Application Programming Interface), which defines important variables such as the kind of requests to be made, how to make them and the data format for transmission. Importantly, this permits communication between

different software medias. An API is entirely customizable, and can be made to extend existing functionality or tailor-made for specific user-demanding modules.

The APIs that will be encountered is handled by the use of *HTTP* (Hypertext Transfer Protocol), which in its simplest form is a protocol that allows the fetching of resources. The protocol is client-server based, such as the client is requesting information and the server is responding to the request.

The most common HTTP-methods are GET, POST and HEAD, which are used to either retrieve, send, or get information about data, respectively. The latter request is usually done before a GET-method for requests considering large amount of data, since this can be a significant variable for the client's bandwidth and load time. Following a request, the server normally responds with one of the status codes in table 5.1.

Table 5.1: Numeric status code for response. The leftmost digit decide the type of response, while the two follow-up digits depends on the implemented API.

Status code	Description
2xx	OK - request was successful
3xx	Resource was redirected
4xx	Request failed due to either unsuccessful authentication or client error.
5xx	Request failed due to server error.

A RESTful (Representational State Transfer) allows users to communicate with a server via a HTTP using a REST Architectural Style [92]. This enables the utilisation of Uniform Resource Identifiers (URI), where each object is represented as a unique resource and can be requested in a uniform manner. Importantly, this allows the use of both URIs and HTTP methods in an API, such that an object is represented by a unique URI whereas a HTTP-method can act on the object. This action will then return either the result of the action, or structured data that represents the object.

To provide a Python example, we can check the response by doing a GET request at the database Materials Project RESTful API in code listing 5.1. We use the preamble to version 2 of Materials Project, and add an API-check and an API-key. The response is shown in code listing 5.2. From the output, it is possible to tell that the supplied API-key is not valid, however, the request is valid.

```

1 import requests
2 preamble = "https://www.materialsproject.org/rest/v2/"
3 url = preamble + "api_check"
```



```

4  params = {"API_KEY": "unique_api_key"}
5  response = requests.get(url=url, params=params)
6  print(response.json())

```

Listing 5.1: Practical example of getting a response from Materials Project database.

```

1  {"valid_response": True,
2  "response":
3    {"api_key_valid": False,
4     "details": "API_KEY is not a valid key.",
5     "version":
6       {"db": "2020_09_08",
7        "pymatgen": "2020.8.13",
8        "rest": "2.0"}
9    }
10 }

```

Listing 5.2: Practical example of response from Materials Project request based on 5.1. The request was done 28. january 2020.

Table 5.2: A selection of databases of computational material science sorted after number of compounds. Abbreviations used are Novel Materials Discovery (NOMAD), Automatic-FLOW for Materials Discovery (AFLOW), Materials Project (MP), Inorganic Crystal Structure Database (ICSD) and Open Quantum Materials Database (OQMD). The number of entries can give the wrong perception of size of each respective database, as it does not visualise how many calculations have been done for each entry, nor if there might be duplicates.

Database	API	Free educational access	Number of entries
AFLOW	REST	True	3.27 M
OQMD [93, 94]	RESTful API (qmpy, matminer)	True	0.82 M
MP [95]	MAPI [96]	True	0.71 M
ICSD [97]	RESTful API	False	0.21 M
Jarvis-DFT	API	True	0.04 M

5.1.2 Open-source Python libraries for material analysis

Many of the databases share convenient modules that are used to adapt, visualize, calculate or predict properties, making it easier for scientists to utilise

Table 5.3: A selection of cloud services that offers database-storage. Abbreviations used are Computational Materials Repository (CMR), NIMS Materials Database (MatNavi), PRedictive Integrated Structural Materials Science (PRISMS), Materials Platform for Data Science (MPDS) and the Materials Data Fascility (MDF).

Cloud service	API/REST	Open educational access
NoMaD	API	True
CMR [98]	ASE	RESTful API
MatNavi	API	True
PRISMS	REST	True
Citrine	API	True
MPDS	API	False
MDF	API	False

the databases. The Atomic Simulation Environment (ASE) is an environment in the Python programming language that includes several tools and modules for setting up, modifying and analyze atomistic simulations [99]. It is in particular used together with the cloud service Computational Materials Repository (CMR).

Another commonly used module is the Python Materials Genomics (pymatgen) [100]. This is a well-documented open module with both introductory and advanced use case examples written in Jupyter Notebook for easy reproducibility, and is integrated with the Materials Project RESTful API.

Another exceedingly popular library is matminer [101], which is an open-source toolkit for material analysis written in Python. Matminer is powered by a group known as *Hacking Materials Research Group*¹. Matminer provides modules to extract data sets from many cloud-services and databases, with examples in table 5.2 and 5.3. Additionally, they provide the tools to extract possibly thousands of features from calculations based on DFT and more, and have modules for visualization and automatic machine learning. These tools will be exemplified in the next chapter.

TODO : Add paragraph about sklearn.

A full selection of python libraries used and their versions can be found in the Github page (TODO: Add github page.)

¹Project's Github site: <https://github.com/hackingmaterials>.

5.2 Databases and cloud services

Every database has its own speciality, and no two databases are the same. There exists entries that are fundamentally identical in several databases, but with different properties as a consequence of parameters used, such as the functional utilised in VASP or the relaxation scheme. This section digs up what exactly is each respective database's claim to fame.

5.2.1 Novel Materials Discovery

The Novel Materials Discovery (NOMAD) [102] Repository is an open-access platform for sharing and utilizing computational materials science data. NOMAD also consists of several branches such as NOMAD Archieve, which is the representation of the NOMAD repository parsified into a code-independent format, NOMAD Encyclopedia, which is a graphical user interface (GUI) for characterizing materials, and lastly NOMAD Analytics Toolkit, which includes early-development examples of artificial-intelligence tools [102].

Databases that are a part of NOMAD data collection includes Materials Project, the Open Quantum Materials Database and AFLOW. They are all based on the underlying quantum engine VASP.

5.2.2 Materials project

Materials project [95] is an open source project that offers a variety of properties of over one hundred thousand of inorganic crystalline materials. It is known as the initiator of materials genomics and has as its mission to accelerate the discovery of new technological materials, with an emphasis on batteries and electrodes, through advanced scientific computic and innovative design.

Every compound has an initial relaxation of cell and lattice parameters performed using a 1000k-point mesh to ensure that all properties calculated are representative of the idealized unit cell for each respective crystal structure. The functional GGA is used to calculate band structures, while for transition metals it is applied +U correction to correct for correlation effects in d- and f-orbital systems that are not addressed by GGA calculations [103]. The thermodynamic stability for each phase with respect to decomposition, is also calculated. This is denoted as E Above Hull, with a value of zero is defined as the most stable phase at a given composition, while larger positive values indicate increased instability.

Each material contains multiple computations for different purposes, resulting in different 'tasks'. The reason behind this is that each computation has a purpose, such as to calculate the band structure or energy. Therefore,

it is possible to receive several tasks for one material which results in more features per material.

5.2.3 AFLOW

The AFLOW[104–106] repository is an automatic software framework for the calculations of a wide range of inorganic material properties. They utilise the GGA-PBE functional within VASP with projector-augmented wavefunction (PAW) potentials to relax twice and optimize the ICSD-sourced structure. They are using a 3000 – 6000 k-point mesh, indicating a more computationally expensive calculation compared to the Materials Project. Next, the band structure is calculated with an even higher k-point density, in addition to the +U correction term for most occupied d- and f-orbital systems, resulting in a standard band gap [107]. Furthermore, they apply a standard fit gathered from a study of DFT-computed versus experimentally measured band gap widths to the initial calculated value, obtaining a fitted band gap [108].

AFLOW-ML [109] is an API that uses machine learning to predict thermomechanical and electronic properties based on the chemical composition and atomic structure alone, which they denote as *fragment descriptors*. They start with applying a classification model to predict if a compound is either a metal or an insulator, where the latter is confirmed with an additional regression model to predict the band gap width. To be able to predict properties on an independent data set, they utilise a fivefold cross validation process for each model. They report a 93% prediction success rate of their initial binary classification model, whereas the majority of the wrongful predictions are narrow-gap semiconductors. The authors does not compare their predicted band gap to experimental values, but it is found that 93% of the machine-learning-derived values are within 25% of the DFT +U-calculated band gap width [110].

5.2.4 Open Quantum Materials Database

The Open Quantum Materials Database (OQDM) [93, 94] is a free and available database of DFT-calculations. It has included thermodynamic and structural properties of more than 600.000 materials, including all unique entries in the Inorganic Crystal Structure Database (ICSD) consisting of less than 34 atoms.

The DFT calculations are performed with the VASP software whereas the electron exchange and correlation are described with the GGA-PBE, while using the PAW potentials. They relax a structure using 4000 – 8000 k-point mesh, indicating an even increasing computational expensive calculation than AFLOW again. Several element-specific settings are included such as using

the +U extension for various transition metals, lanthanides and actinides. In addition, any calculation containing 3d or actinide elements are spin-polarized with a ferromagnetic alignment of spins to capture possible magnetism. However, the authors note that this approach does not capture complex magnetic, such as antiferromagnetism, which has been found to result in substantial errors for the formation energy [111].

5.2.5 JARVIS

Joint Automated Repository for Various Integrated Simulations (JARVIS) [112] - DFT is an open database based on the VASP software to perform a variety of material property calculations. It consists of roughly 40,000 3D and 1,000 2D materials using the vdW-DF-OptB88 van der Waals functional, which was originally designed to improve the approximation of properties of two-dimensional van der Waals materials, but has also shown to be effective for bulk materials [113, 114]. The functional has shown accurate predictions for lattice-parameters and energetics for both vdW and non-vdW bonded materials [115].

Structures included in the data set are originally taken from the materials project, and then re-optimized using the OPT-functional. Finally, the combination of the OPT and modified Becke-Johnson (mBJ) functionals are used to obtain a representative band gap of each structure, since both have shown unprecedented accuracy in the calculation of band gap compared to any other DFT-based calculation methods [116].

The JARVIS-DFT database is part of a bigger platform that includes JARVIS-FF, which is the evaluation of classical forcefield with respect to DFT-data, and JARVIS-ML, which consists of 25 machine learning to predict properties of materials. In addition, JARVIS-DFT also includes a data set of 1D-nanowire and 0D-molecular materials, yet not publically distributed.

Chapter 6

Information flow

The information stream of this project can be regarded as many modular parts connected together in logical pieces, and is strongly influenced by the process that defines a *minimum viable product* (MVP) through iterative development. An MVP is commonly known (in the business world) as a new product that enables the most learning out of the minimum effort possible. This method allows a product to be iteratively evolved by consistent feedback and development, which in return enables cooperation between cross-disciplinary fields.

Furthermore, by having several modules serving as the fundament of the project, it is possible to achieve a long-lasting and robust product that is simple to maintain yet straightforward to develop. Bugs can be tackled through a documented code simultaneously as visible future improvements can be addressed. Therefore, the product is not regarded as completed in any terms, but rather ready for a first release after iteratively finding the minimum viable product.

The main project of this work can be found on the Github repository *predicting-solid-state-qubit-candidates* [117]. In this chapter we will look into the details and thoughts behind the extraction of data, building features, data preparation, data mining and eventually fabricating a generalized model that can predict unseen data with confidence.

6.1 Extraction and featurization of data

The initial step for gathering and building features can be visualised through the flowchart in figure 6.1. Initially, we start by extracting all entries in the Materials Project that matches a specific query. Thereafter, we apply Matminer's featurization tools to make thousands of features of the data. In a parallel step, entries that are deemed similar to the entries from the initial Materials Project query are extracted from AFLOW, AFLOW-ML, JARVIS-DFT, OQMD and Citrine Informatics. Finally, we combine the steps together

as interim data that is ready for further analysis.

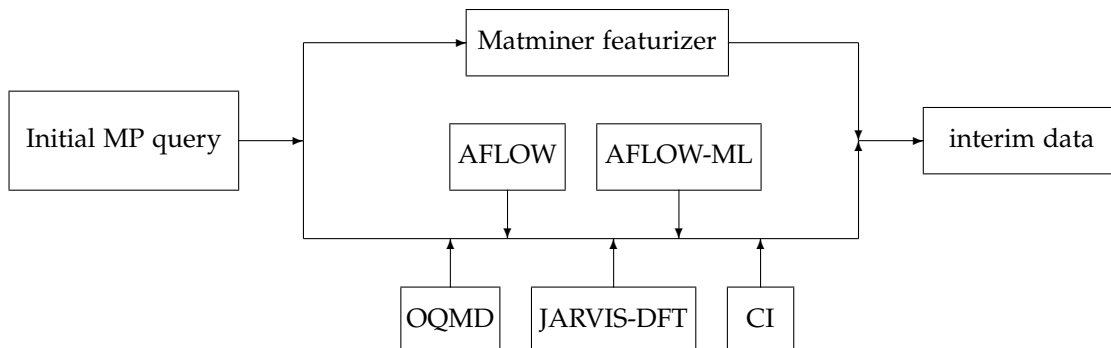


Figure 6.1: The data flow of the main project, starting from an initial MP-query, and ending with a featurized dataset with entries from several other databases. The matminer featurizer step is further visualized in-depth in figure 6.2.

The initial query has the requirement that all entries has to be derived from an experimental ICSD entry, and is reasoned by that we can identify equivalent entries in other databases. Furthermore, all entries in the Materials Project needs to have a band gap larger than 0.1eV. Recall that Materials Project applies the functional GGA in estimating the band gap, which is known to severely underestimate the given electronic property. Therefore, we have chosen a low value to not rule out any potential candidates but high enough to leave out all materials that can be considered metallic. Thus, out of a total of 139.367 entries in Materials Project, our initial requirement is satisfied by 25.352 of the entries.

From figure 6.1 we notice that by using many databases we do not add additional entries that exist in some databases but is not to be found in Materials Project. This is by design since it preserves the versatility of choosing a database to work with. Therefore, one can completely ignore steps such as the initial query of Materials Project or the featurization process, and rather focus on e.g. all the 400.000 entries existing in OQMD. The examples that follows will illustrate the ease of extracting data from several different databases, and can serve as the starting point for other research projects in computational material science.

6.1.1 Practical data extraction with Python-examples

For this section, we will show practical examples of how to extract data that might fulfill the criteria for a material to host a qubit candidate given in the theory part. We will begin with the database of Materials Project, and then

search for entries in other databases that match entries from MP. This process is reproducible as a jupyter notebook¹ and the databases in question are the ones referred to in the previous section.

Instead of building multiple HTTP-methods from scratch, we will here take a look at the easiest method at obtaining data from each database. The range of data in a database can consist of data from a few entries up to an unlimited amount of entries with even further optional parameters, and has limitless use in applications. However, the amount of data in a database is irrelevant if the data is inaccessible. Therefore, we provide a toolbox in how to extract information in the easiest way possible. This includes looking into the APIs that supports data-extraction and that are recommended by each respective database.

Every data extraction class is based on an abstract parent class. The advantages of using a base parent class are many, such as improving the readability during code reviews, reducing the main barrier for understanding the underlying structure of a project and utilising reusable components. Yet, the main advantage of using a base parent class is the fact that it can effortlessly be extended for further implementations since it provides a code skeleton.

Materials Project

The most up-to-date version of Materials Project can be extracted using the python package pymatgen, which is integrated with Materials Project REST API. Other retrieval tools that is dependent on pymatgen includes matminer, with the added functionality of returning a pandas dataframe. Copies of Materials Project are added frequently to cloud services such as Citrine Informatics, but the latest added entries to Materials Project cannot be guaranteed in such a query.

Entries in Materials Project are characterized using more than 60 features², some features being irrelevant for some materials while fundamental for others. The data is divided into three different branches, where the first can be described as basic properties of materials including over 30 features, while the second branch describes experimental thermochemical information. The last branch yields information about a particular calculation, in particular information that's relevant for running a DFT script.

To extract information from the database, we will be utilising the module pymatgen. This query supports MongoDB query and projection operators³, resulting in an almost instant query.

¹add and insert DOI for JN 01-generateDataset-notebook.ipynb

²All features can be viewed in the documentation of the project: <https://github.com/materialsproject/mapidoc/master/materials>

³<https://docs.mongodb.com/manual/reference/operator/query/>

1. Register for an account⁴, and generate a secret API-key.
2. Set the required criteria.
3. Set the wanted properties.
4. Apply the query.

The code nippet in code listing 6.1 resembles steps 2 – 4, and is filtered as the initial query.

```
1 from src.data.get_data_MP import data_MP
2
3 MAPI_KEY = 'very_secret_key_here'
4 MP = data_MP(API_KEY=MAPI_KEY)
5 df = MP.get_dataframe()
```

Listing 6.1: Practical example of extracting information from Materials Project using pymatgen, resulting in a Pandas DataFrame named entries that contains the properties given after performing a filter on the database. The criteria is given as a JSON, and supports MongoDB operators.

Citrine Informatics

Citrine Informatics is a cloud service, which means that the spectrum of stored information varies broadly. We will access research through open access for institutional and educational purposes. Information in Citrine can be stored using a scheme that is broken down into two sections, with private properties for each entry in addition to common fields that are the same for all entries.

In this example, we will gather experimental data using the module matminer. The following steps are required to extract information from Citrine Informatics.

1. Register for an account⁵, and generate a secret API-key.
2. Set the required criteria.
3. Set the wanted properties and common fields.
4. Apply the query.

The code listed in code listing 6.2 gives an easy example to steps 2 – 4 with experimental data as filter, which results in an almost instant query.

⁴<https://materialsproject.org>

⁵<https://citrination.com>

```
1 from src.data.get_data_Citrine import data_Citrine
2
3 CAPI_KEY = 'very_secret_key_here'
4 citrine = data_Citrine(API_KEY=CAPI_KEY)
5 df = citrine.get_dataframe()
```

Listing 6.2: Practical example of extracting information from Citrine Informatics using matminer, resulting in a Pandas DataFrame named `experimental_entries` that contains the properties given after performing a filter on the database. The criteria is given as a JSON.

AFLOW

The query from AFLOW API [104] supports lazy formatting, which means that the query is just a search and does not return values but rather an object. This object is then used in the query when asking for values. For every object it is necessary to request the desired property, consequently making the query process significantly more time-demanding than similar queries using APIs such as pymatgen or matminer for Citrine Informatics. Hence, the accessibility is strictly limited to either searching for single compounds or if the user possess sufficient time.

Matminer's data retrieval tool for AFLOW is currently an ongoing issue [118], thus we present in code listing 6.3 a function that extracts information from AFLOW and returns a Pandas DataFrame. In contrast to Materials Project and Citrine Informatics, AFLOW does not require an API-key for a query, which reduces the amount of steps to obtain data. The class searches for an stored AFLOW-data, and initialises a MP-query with the initial criteria if not successful. The resulting query will then be used as input to AFLOW.

```
1 from src.data.get_data_AFLOW import data_AFLOW
2
3 AFLOW = data_AFLOW()
4 df = AFLOW.get_dataframe()
```

Listing 6.3: Practical example of extracting information from AFLOW. The function can extract all information in AFLOW for a given list of compounds, however, it is a slow method and requires consistent internet connection.

Restricted by the available API, the resulting query of 25212 entries in Materials Project took place during the period from january to february 2021 and took in total 23 days. Unfortunately, less than 0.02% of the entries screened from Materials Project was present in AFLOW.

AFLOW-ML

In this part, we will be using a machine learning algorithm named AFLOW-ML Property Labeled Material Fragments (PLMF) [109] to predict the band gap of structures. This algorithm is compatible with a POSCAR of a compound, which can be generated by the CIF (Crystallographic Information File) that describes a crystal's generic structure. It is possible to download a structure as a poscar by using Materials Project front-end API, but is a cumbersome process to do so individually if the task includes many structures. Extracting the feature of POSCAR is yet to be implemented in the RESful API of pymatgen, thus we demonstrate the versatility of pymatgen with a workaround.

We begin with extracting the desired compounds formula, their Materials Project IDs (MPIDs) for identification, and their respectful structure in CIF-format from Materials Project. In an iterative process, each CIF-structure is parsed to a pymatgen structure, where pymatgen can read and convert the structure to a POSCAR stored as a Python dictionary. Finally, we can use the POSCAR as input to AFLOW-ML, which will return the predicted band gap of the structure. This iteratively process parsing and converting, but is an undemanding process. The function that handles this is presented in code listing 6.4. Similar to AFLOW-query, this code listing is dependent on MP-data and will apply for a query if the data is not present.

A significant portion of the process is tied up to obtaining the input-file for AFLOW-ML, and fewer structures will result in an easier process. Nevertheless, we present the following steps in order to receive data from AFLOW-ML.

1. Download AFLOWmlAPI⁶.
2. Getting POSCAR from MP.
 - (a) Apply the query from Materials Project with "CIF", "material_id" and "full_formula" as properties.
 - (b) Insert resulting DataFrame into function defined in code listing 6.4.
3. Insert POSCAR to AFLOW-ML.

```
1 from src.data.get_data_AFLOWML import data_AFLOWML
2
3 AFLOWML = data_AFLOWML()
4 df = AFLOWML.get_dataframe()
```

⁶<http://aflow.org/src/aflow-ml/> to the same directory as code listing 6.4

Listing 6.4: Practical example of extracting information from AFLOW-ML. The function will convert a CIF-file (from e.g. Materials Project) to a POSCAR, and will use it as input to AFLOW-ML. In return, one will get the structure's predicted band gap. It should be noted that this requires the AFLOW-ML library in the same directory.

The resulting ab-initio calculations used an average of 57s/compound, which in total sums up to 16.6 days. In contrast to AFLOW, 100% of the entries was present due to the fact that it is not based on a database but rather a model.

OQMD

To extract information from the OQMD, the easiest way was through the interface of Matminer. The difficulty of extraction are mostly regarded to column which are not assigned to a type, however, this is taken care of in the extraction class visualized in code listing 6.5.

```
1 from src.data.get_data_OQMD import data_OQMD
2
3 OQMD = data_OQMD()
4 df = OQMD.get_dataframe()
```

Listing 6.5: Practical example of extracting information from OQMD through Matminer.

The query is done almost instantly, resulting in a DataFrame containing over 400.000 entries, where 40% of the entries are matching an entry of the initial MP query.

JARVIS-DFT

The newest version of the JARVIS-DFT dataset can be obtained by requesting an account at the official webpage, but with the drawback that an administrator has to either accept or deny the request. Thus, the accessibility of the database is dependent on if there is an active administrator paying attention to the requests, which is a limitation experienced during this work. Another approach is to download the database through matminer, however with the limitation of not necessarily having the latest version of the database. A third approach is to download a version of JARVIS-DFT that have been made available for requests the 30.04.2020 at <http://figshare.com> by Choudhary *et al.* [112]. The author provides tools for extraction, yet not compatible with the latest version of Python (3.8) at the time writing (12.03.2021). Therefore, we provide a tool to extract this data through the use of our base class.

```
1 from src.data.get_data_JARVIS import data_JARVIS
2
3 JARVIS = data_JARVIS()
4 df = JARVIS.get_dataframe()
```

Listing 6.6: Practical example of extracting information from JARVIS-DFT. For this example, we exclude all metals by removing all non-measured band gaps.

We observe that there is no advanced search filter when loading the database from matminer. The author of matminer regards this as the user’s task, and is indeed easily done through the use of the python library Pandas.

The resulting screening of 25212 entries from Materials Project was done almost instantly, and it was found 11% and 17.8% similar entries for the TBMBJ and OptB88 functionals with MP, respectively. Moreover, JARVIS-DFT contains information about spin-orbit splitting, but only 0.12% of the calculations was found as a match with the initial MP query.

6.2 Matminer featurization

Before applying any machine learning algorithm, raw data needs to be transformed into a numerical representation that reflects the relationship between the input and output data. This transformation is known as generating descriptors or features, however, we will in this work adapt the name *featurization*. The open source library of Matminer provides many tools to featurize existing features extracted from Materials Project. In this section we will describe how to extract the features from an initial Materials Project query result (see subsection. 6.1.1), and the resulting features. It is beyond the scope of this work to go in-depth of each feature since the resulting dataset contains a quantity of more than 4500 features, but we will here take the liberty to serve a brief overview of the features and refer to each respective citation for more information. The respective table with information regarding 39 distinct matminer featurizers is situated in the Appendix, table B.1.

The motivation behind the choice of featurizers is that we do not precisely know which features that describes a good potential host. A few potential candidates were briefly mentioned in section 2.6.4, while most candidates are probably yet to be discovered. If we had precise knowledge of what to look for, then there is a good chance that the list of hosts would be longer. Therefore, we strive to collect an achievable quantity of descriptors with the hope of getting wiser in terms of describing a potential material host.

To apply matminer’s featurization tools, we extend an existing implementation by Breuck *et al.* [119] called the Materials Optimal Descriptor Network

(MODNet). The author Breuck *et al.* specifies that MODNet is a supervised machine learning framework for learning material properties based on either composition or crystal structure. To provide the training data for their model, MODNet featurizes (through matminer) structures either from Materials Project or in the form of a structure object made by pymatgen. Their current implementation provides featurization for compositions, structures and sites. However, matminer also provides featurization tools for density of states (DOS) and band structures, therefore we modify MODNet and extend it to facilitate such featurizations.

One immediate limitation of our extension is that Matminer's tools is dependent on a pymatgen DOS- and bandstructure object. These objects contains information up to 10MB, and becomes a challenge when dealing with data containing several thousand such objects. This is solved by the required features for matminer's featurization for a subsample of the data, followed by a featurization process of the same subsample. When the featurization is done, we store the new features and throw away the pymatgen features. This is done iteratively for the entire data set. Thus, a compromise between applying several queries and storing information has been done. The scheme can be visualised as the flow chart seen in figure 6.2.

In the extended version of the featurization process, we eliminate all columns that does not have any entries with physical meaning. This is beneficial for several reasons, such as to reduce memory allocated and to preprocess the data. If there are entries existing with both physical and non-physical for the same column, we replace the non-physical meanings with -1 for recognition in a later step. Additionally, we convert columns that are categorical or lacks a numerical representation into a categorical portrayal. Thus, we strive to limit the necessary steps for further processing of data into a machine learning algorithm. Nevertheless, the featurization process results in 4876 descriptors.

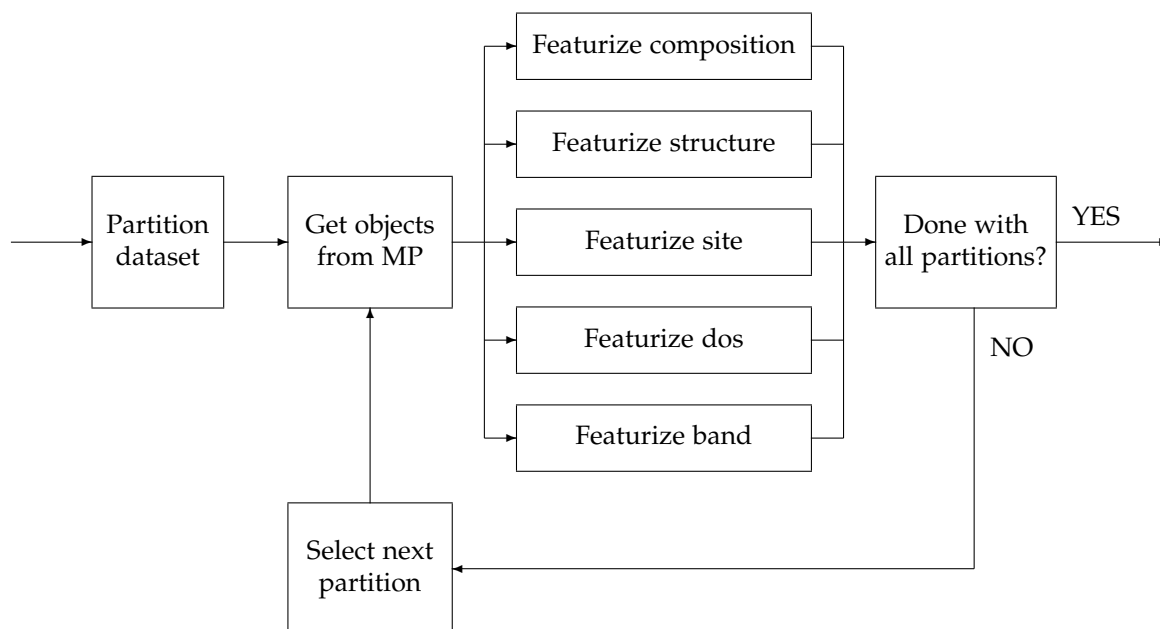


Figure 6.2: The process of the matminer featurizer step as seen in figure 6.1. To limit the memory and computational usage, the data is partitioned into smaller subsets where the respective pymatgen objects are obtained through a query to be used in the following featurization steps. This is iteratively done until all the data has been featurized.

Even if the first version of Matminer was released in 2016, many issues concerning daily operational use are still present. During the featurization process in this work, we manually identified 14 (TODO: Update number) erroneous entries that are summarized in the Appendix, table B.2, which were excluded from the dataset. These entries were part of the reason why the featurization process is a time-consuming process, as there is currently no implementation in Matminer that can potential pick up and catch erroneous entries in Materials Project. The process of manually catching such an entry was identified by featurization of single entries causing one of two problems. The first problem could be that an entry could be causing a memory leak which leads to an exceedingly large memory allocation, or it could be that the featurization process needed days to calculate oxidation states for a structure.

6.3 Labelling of data

After selecting entries based on an initial query from Materials Project followed by a thorough featurization process using Matminer, we face a challenge in terms of defining a training set that we can train data on. This is not only challenging due to the lack of known candidates, but also due to the intricacy of defining materials as bad candidates. Therefore, in this section we describe three different approaches of finding a training set consisting of (1) good candidates and (0) bad candidates.

6.3.1 First approach; the Ferrenti approach

The first approach on defining a training set is based on the criteria from the paper “Identifying candidate hosts for quantum defects via data mining” of Ferrenti *et al.* [110], therefore we will name this approach *the Ferrenti approach*. They suggest a data mining process consisting of four stages by systematically evaluating the suitability of host materials from Materials Project. This procedure is referred to as *data mining*, and we will initially begin with looking at labelling good candidates.

Labelling good candidates

The first stage consists of the following steps to include materials that

Stage 1

- contains elements with a $> 50\%$ natural abundance of zero spin isotopes.
- crystallize in nonpolar space groups.
- is present in the ICSD database.
- is calculated nonmagnetic.

The restriction of materials to only contain elements with at least 50% nuclear spin-free isotopes might help with reducing decoherence for all semiconductor-based quantum technologies, as discussed in section 2.6.2. The limit is chosen due to that elemental species with at least 50% nuclear spin free isotopes could likely be isotopically enriched to higher concentrations [110], which has been accomplished for carbon [120, 121] and silicon [122]. In particular, the restriction excludes the use of 53 elements from any species. Any magnetic noise or any presence of electric dipole moment could also potentially increase decoherence of defects. Therefore, we try to reduce any decoherence by restricting materials to possess highly symmetric structures which are nonmagnetic.

Stage two consists applying additional filtering due to practical reasons. This includes removing all materials containing radioactive or toxic elements, as well as removing noble gases because none exists as solids under standard conditions. Rare-earth metals were also excluded due to the difficulty of obtaining pure materials that are sufficiently free of nuclear spin. Lastly, we remove entries that occur mostly in very complex cluster structures (Ru, Os) or are not present in any identified phases (Fe, Ni). Therefore, the additional filter constitutes of obtaining materials that

Stage 2

- does not include Th, U, Cd or Hg.
- does not include any noble gases or rare-earth elements.
- does not include Ru, Os, Fe, Ni

Stage three consists of setting a lower band gap limit similar to that of silicon, but due to severe underestimation of bandgaps by PBE-GGA we set this restriction lower since we do not want to exclude any potential host candidates. The materials are required to have

Stage 3

- a bandgap larger than 0.5eV as calculated by MP PBE-GGA.

Finally, the last stage consists of identifying the thermodynamic stability of each compound. A large energy above hull (E Above Hull) is an indication of an unstable compound and would likely cause decomposition, therefore the last filter requires the materials to have

Stage 4

- a calculated E Above Hull $< 0.2\text{eV/atom}$.

The quantity of entries through the different stages have been visualized in table 6.1. The table compares our and their implementation of the same screen procedure with different results. In particular, we see that the remaining materials that have survived four stages of filtering are twice as many. This could be due to the date of extracting since it differs with 13 months, since over 14.000 new entries were added to Materials Project. However, another reason could be due to that they have done additional manual screening. Unfortunately, precise information of which entries that were excluded from the manual filtering were not included in neither the article or the supplementary information [110]. Yet, after doing a data mining procedure we have found 1046 potential candidates that exhibit promising features.

Table 6.1: A table that compares two different implementations of the same screen procedure. Ferrenti *et al.* extracted information March of 2020, while we did the extraction during April of 2021. The adjusted difference is given as our reported entries divided on their reported entries.

Stage	Good candidates based on Ferrenti <i>et al.</i> [110]	Good candidates based on approach 1	Adjusted difference
Total entries in Materials Project [95, 96]	125.223	139.367	11%
Stage 1	3363	4347	29%
Stage 2	1993	2226	12%
Stage 3	920	1181	28%
Stage 4	541	1046	93%

Labelling bad candidates

We have now defined good candidates, and turn our attention to defining bad candidates. This is perhaps the difficult part, since we do not know exactly which properties or combination of features a material needs to exhibit for it to be excluded from any use in quantum technology. Therefore, we try to find the opposite criteria of the four stages that defined good candidates.

If we were to turn around all criteria defined in the four stages above (except for energy above hull), it would result in only 52 entries which would make the combined data set very imbalanced. Instead, we try to provide a more general process that includes a larger variety of entries, which could potentially increase the predictor space for bad candidates. The screening procedure for stage 1 requires bad candidates to

Stage 1

- crystallize in polar space groups.
- be present in the ICSD database.
- be calculated as magnetic.

Stage 2

- have a bandgap larger than 0.1eV as calculated by MP PBE-GGA.

We include only ICSD entries and a lower band gap limit for consistency, since our data does not contain entries outside of these limits. The number of entries after stage 1 is 1520, while stage 2 reduces the entries to 684.

6.3.2 Second approach; the augmented Ferrenti approach

In the second approach we try to make adjustment of the first approach to improve the dataset. This approach is therefore named *the augmented Ferrenti Approach*.

Labelling good candidates

The first approach included unphysical criteria such as removing elements that are either radioactive, toxic, elements not occurring under standard conditions, or rare-earth elements that are difficult to obtain. In this approach, we remove those constraints since these are not criteria that necessarily deem a material as either good or bad for QT, and it is eventually up to experimentalists for evaluation of such practicalities. Therefore, we remove stage 2.

We will in this approach not consider if a material is stable or not, since this is eventually up to experimentalists to evaluate. Additionally, we will include a few interesting elements that showed promising properties as discussed in section 2.6.4, and was originally excluded due to lack of spin zero isotopes. Thus, the second approach consists of the following steps to include materials that

Stage 1

- contains elements with a $> 50\%$ natural abundance of zero spin isotopes except Al, P, Ga, As, B and N.
- crystallize in nonpolar space groups.
- is present in the ICSD database.
- is calculated nonmagnetic.

Stage 2

- have a bandgap larger than 1.5eV as calculated by MP PBE-GGA.

Stage 3

- have a calculated E Above Hull $< 0.2\text{eV/atom}$.

Since we have removed restrictions, we can also infer a stronger one for the band gap. Therefore, we can be considerable more certain if a band gap can accommodate deep defect due to an increasing amount of entries when removing restrictions.

Labelling bad candidates

For bad candidates, we implement the same strategy as defined for bad candidates in approach 1. The resulting table for both good and bad candidates is found in table 6.2. The table reveals a considerable imbalanced dataset with up to 75% being good candidates, while only 25% of the training data are labelled as bad candidates. However, the training set is 78% larger than in approach 1.

Table 6.2: A table showing the number of entries through the data mining process for good candidates in approach 2 and bad candidates in approach 1 and 2.

Stage	Good candidates approach 2	Bad candidates approach 1 and 2	Ratio
Total entries in Materials Project [95, 96]	139.367	139.367	-
Stage 1	7433	1520	83%/17%
Stage 2	2373	684	78%/22%
Stage 3	2141	—	75%/25%

6.3.3 Third approach; the insightful approach

The third approach is vastly different than the two first approaches in terms of labelling, therefore it is named *the insightful approach*.

Recall, in section 2.6.4 we discussed alternative promising material host candidates. The third approach for finding good candidates is to search our current data for any materials that overlap with known good candidates. Due to a concern of having a too small dataset, we will include materials that are promising and have shown suitable properties to accomodate deep defects that can exhibit quantum effects.

Labelling good candidates

Stage 1

- matches the formulas SiC [6, 30, 37, 46, 47], BN [56, 57], MoS₂[57], WSe₂[57], WS₂[57], GaN [52], GaAs [51], AlN [6, 53], ZnS [49], ZnSe [6], ZnO [49], AlP[6], GaP[6], AlAs[6], ZnTe[6], CdS[6], SiGe [55], C [33, 35, 36] or Si [49, 50].
- is present in the ICSD database.

Stage 2

- have a bandgap larger than 0.5eV as calculated by MP PBE-GGA.

Stage 3

- Manual screening of correct structures.

After stage 1, it was found 202 matching formulas which included 12 entries that had a bandgap lower than 0.5eV. These entries were structures that was reported as unstable in terms of energy above hull calculations, and would decompose into entries that were already present in the data after stage 1 with bandgap substantially larger than 0.5. We choose to exclude these entries with an additional band gap restriction due to the fact that the bandgap is not large enough to accomodate a deep defect. Therefore, these entries were instead labelled as bad entries.

Entries matching the formula C, SiC, BN, MoS₂, WSe₂ and WS₂ were manually screened to see if the entries have a matching structure to the respective candidates discussed in section 2.6.1 and 2.6.3 and 2.6.4, respectively. For C, we admit only three-dimensional diamond-like structures, as explicitly stated in the column tags at Materials Project. Two-dimensional graphite-like structures are labelled as bad candidates, while complex structures (eg. C₂₈, C₄₈, C₆₀) were moved to the test set. For SiC, we admit only entries with the polytypes 3C, 4H and 6H, while moving structures similar to 2H to the test set. Concerning the other materials, we only admit two-dimensional structures.

The materials AlP, GaP, AlAs, ZnTe and CdS were manually screened for tetrahedrally coordinated structures, and have been included since Weber *et al.* [6] has identified them as potential promising candidates due to acceptable properties defined in requirements (H1-H4) in section 2.6.2. We note that only tetrahedrally coordinated structures of the given formulas were present after the bandgap restriction of 0.5eV.

Since the number of elements in the good candidates are not containing more than two elements, we decide to remove the feature that explains how many elements due to that we do not want the model to discriminate based on this feature. After three stages, a total of 172 entries were labelled as good candidates.

Labelling bad entries

Since the training data that constitutes good candidates are few, we choose to add 400 random entries from the dataset of bad candidates used in approach 1 and 2, in addition to the entries stated above. We only add a subsample for increasing the potential dimensional space for predictions of candidates while avoiding having a too inbalanced dataset.

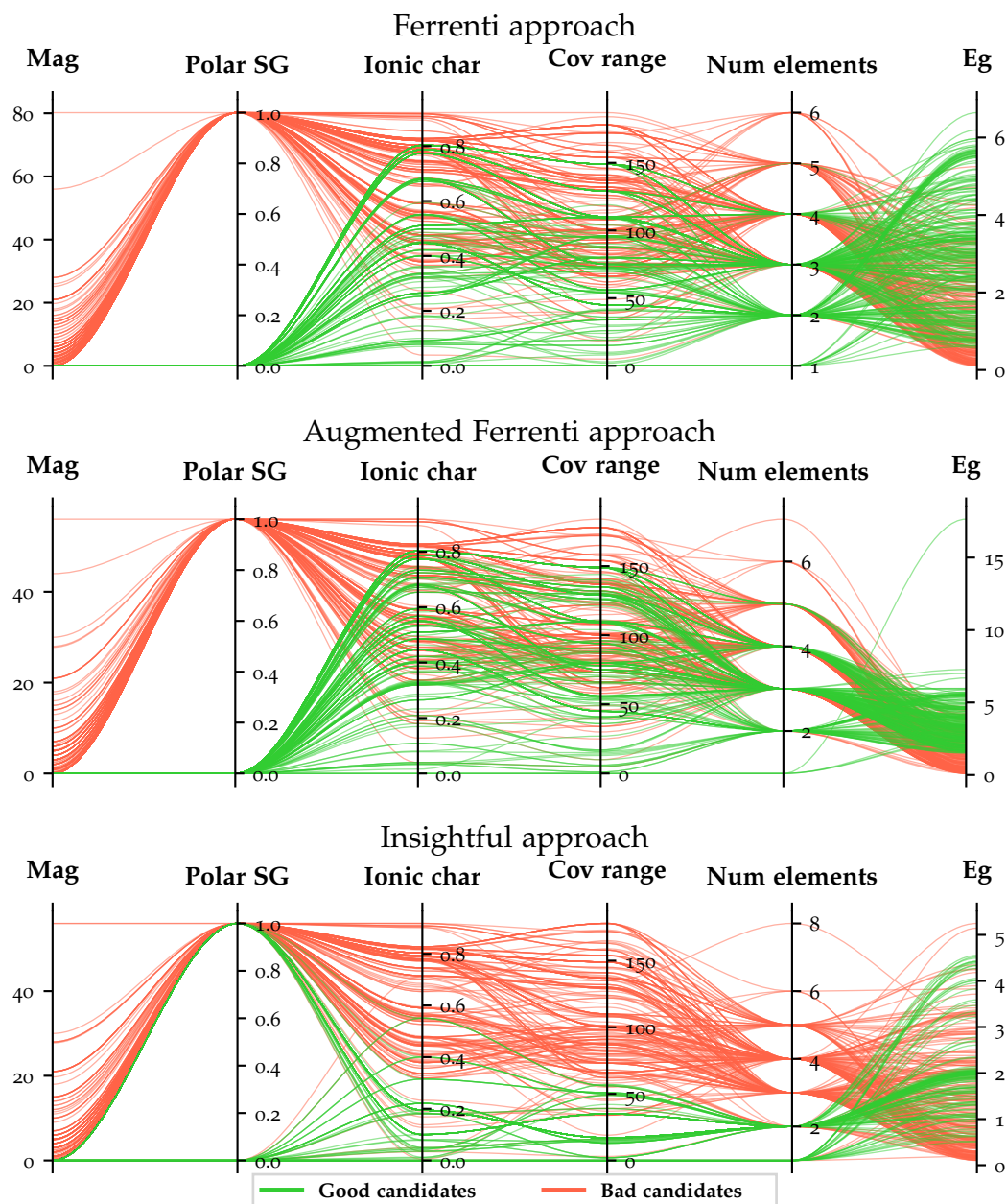


Figure 6.3: Parallel coordinate plots for the different approaches. To limit the data cluttering, we have random collected up to 250 entries for each class and made the lines transparent. For the insightful approach, we have used all entries 172 entries.

6.3.4 Comparison of the approaches

The three approaches provide special emphasis on each their goal. The Ferrenti approach is dependent on choosing only elements with zero spin isotopes together with practical filters, while the augmented Ferrenti approach allows a larger variety of elements and removes the practical reasons for excluding elements. Thus, the first approach targets a more narrow prediction space than the second approach does, and we would expect that the second approach will lead to more predicted candidates compared to approach one. However, perhaps the most restricted approach is the insightful approach. Since we only include known candidates, they should share the same properties and therefore provide a very narrow prediction space.

Unfortunately, the downside of including all the known candidates in one approach is that it becomes increasingly challenging to evaluate the approach or the resulting model. For the two first approaches, we can see if some of the known candidates are present in the predictions, while this is not possible for the latter approach.

We provide a visualization of each approach's training data as a parallel coordinate plot for a few selected features in 6.3. Parallel coordinate schemes [123, 124] represents a multi-dimensional data tuple as one polyline crossing parallel axis. The selected features are found on the x-axis, while the y-axis show the value of the data present. Thus, parallel coordinate plots can turn complex many dimensional data into a compact two-dimensional representation. However, due to data cluttering and that one entry can potentially reserve a large visual area of the figure, the utilization becomes limited when facing large datasets [125]. Therefore, we have chosen to plot a random sample of each class with an upper limit of 250 per class with transparent lines.

The Ferrenti approach and the augmented Ferrenti approach share similarities, such as only having polar space groups present and having an equal amount of upper limit for both ionic character and covalent range. Additionally, they share that entries constitute of up to five different elements. Interestingly, we can see that even if the augmented Ferrenti approach is less restricted, it appears that the entries map over the same dimension based on 6.3.

The biggest difference is seen for the insightful approach. The chosen entries do not possess any magnetization, even if there are both polar and nonpolar space groups present. The range of covalent radius and maximum ionic character is significantly lower than the two other approaches.

To visualize the complexity of the training sets, we have found the two largest eigenvalues of the covariance matrix of the initial data from Materials Project, and transformed the training sets according to the corresponding two eigenvectors. The resulting scatter plots is found in figure 6.4. In green squares, we find the good candidates for each approach, while the labelled

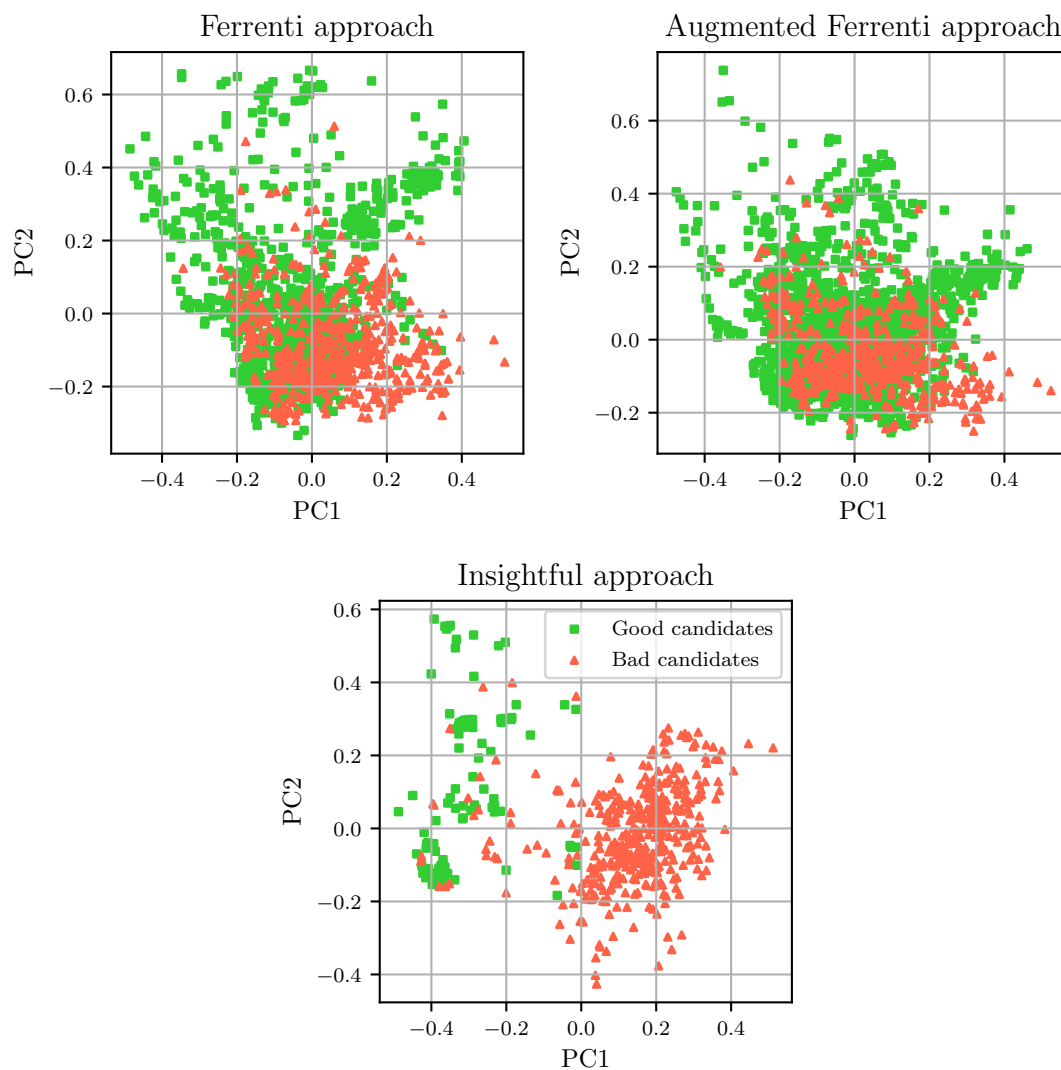


Figure 6.4: Two-dimensional scatter plots for the three different approaches. We have found the two eigenvectors corresponding to the two largest eigenvalues of the covariance-matrix, that is the two most important principal components, of the initial data from the Materials Project query. Then, we have transformed the three training sets resulting from the three approaches and visualized it as a scatter plot for visualization purposes.

bad candidates are dressed in red triangles. Due to the simplicity of reducing the number of features down to 2 features, both good and bad candidates for the Ferrenti approach are overlapping which could be challenging for any model that would try to learn a clear-cut boundary. However, for the insightful approach, we can already start to see a trend where the upper left part of the figure is dominated by good candidates. Therefore, we can expect that the two Ferrenti approaches would need either supplementary dimensions for further distinguishment, or could be in trouble of finding a generalized model.

6.4 Model selection

After building a dataset through extraction, featurization and labelling, we turn our attention towards training a model. The flowchart is visualized in figure 6.5. After gathering and featurization, we achieve the interim data that goes through a final data preparation step to become preprocessed data. Then, we perform a data mining step using three different approaches as discussed in the last section. For each of the three approaches, we train and predict in the step called supervised learning. In the summary, we compare the different approaches and results.

In the data preparation step, we assess the quality of the data. Due to the large dimension of 25000×4500 , we assume that there is a large amount of non-physical values present, therefore we fill all the missing values with zero and remove all columns with more than 70% containing only zeros. This value was chosen since all categorical features have at least 30 of an respective class present in the column, and a majority of the removed columns contained between 90% and 100% only zeros. This reduces the dimensionality substantially to only 679 features. It should be noted that other methods of data preparation resulted in equivalent preprocessed data due to the large amount of missing values in the data.

Four different supervised models has been selected for each of the three approaches defined in the previous section, resulting in a total of 12 unique models. As discussed in section 4.2, models are unique and does not necessarily perform optimal on all kinds of data. Therefore, the four models have been selected as a function of increasing complexity and ranges from the simplistic logistic regression and decision trees and up to random forest and gradient boost. We utilize the implementation of sklearn for all models [126].

Due to the fact that the current dimension of the entire dataset is still large, we apply the dimensionality reduction technique PCA to the dataset. This is beneficial for several reasons, such as finding correlated features and reducing any dimensional. Additionally, it opens up for a visualized interpretation if

we were to choose 3 or less principal components.

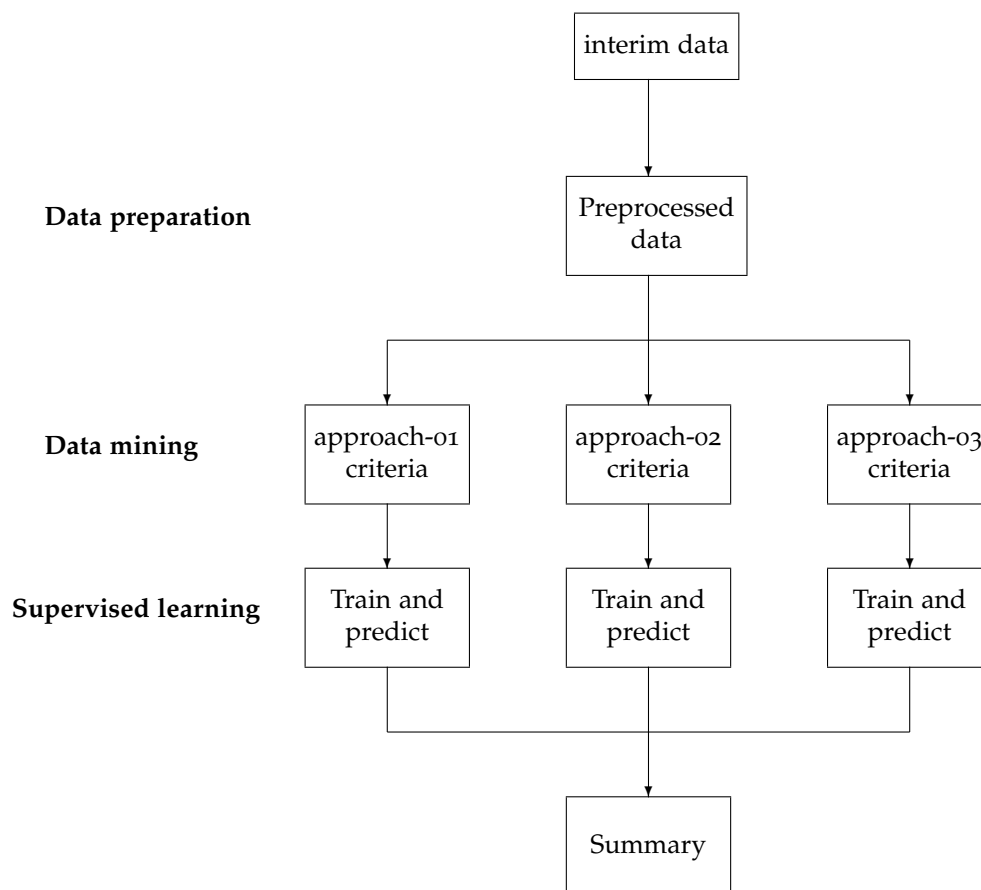


Figure 6.5: A continuity of the flowchart in figure 6.1 that visualise the steps of data preparation, the three approaches in the data mining step and the subsequent supervised learning. Finally, a summary will be provided. From a hierarchic perspective, we find the steps leading up to the preprocessed data as the top level, while each of the approaches are found one level down.

The optimal parameters are then searched for with the use of sklearn's gridsearch and imblearn's pipeline [126, 127]. Imblearn's pipeline enable the use of resampling methods, in contrast to sklearn's pipeline, but does not differ in any other way. In the pipeline, we provide a standardscaler that scales the data such that every feature will have a mean of 0 and a standard deviation of 1 [126]. Thereafter follows the dimensionality reduction and a supervised learning algorithm. It should be noted that the number of optimal principal components, up to an upper limit of accumulated explained variance of 95%, are also searched for in the grid-search scheme.

Part IV

Results and discussion

Chapter 7

Optimization

This chapter is named optimization due to its contents; here we will account for the choices we compose to optimize a potential machine learning algorithm. Initially, that involves finding what information is stored within the databases and the compromise of gathering the information, which further evolves into finding optimal hyperparameters for each approach.

The first step of this work was to find data from Materials Project, involving entries that are associated with an ICSD structure and have a PBE-GGA calculated band gap of minimum 0.1eV. Out of 126.335 existing entries in Materials Project, 48.644 (39%) were found to have an associated ICSD-structure, while 65.783 (52%) materials had a calculated band gap of at least 0.1eV. It was found that 25271 (20%) materials have the band gap minimum and an associated ICSD-structure. It should be noted that these numbers are based on data extraction in December of 2020, while the extraction from other databases and featurization related to this work was done in the time period of December 2020 to March 2021. In February of 2021, over 30.000 new materials were added in a large update¹. These new entries are not included in this work, and therefore the number of entries in our Materials Project are based on the latest release in 2020, which is named V2020.09.08.

Two visualizations of two different distributions of the data is found in figure 7.1 and 7.2. The first figure visualize the distribution of oxid types as a function of compound type, and reveal that the majority of compounds are either binary, ternary, quaternary or quinary, where the majority of the materials are either oxide or not. This is important to know considering our labelling approaches, in particular the insightful approach where we hand-picked good entries. None of the entries labelled as good candidates were oxides, which means that potentially half of the the test set could potentially be considering as bad candidates by the models.

¹<https://matsci.org/t/materials-project-database-release-log/1609/16> 23.04.2021

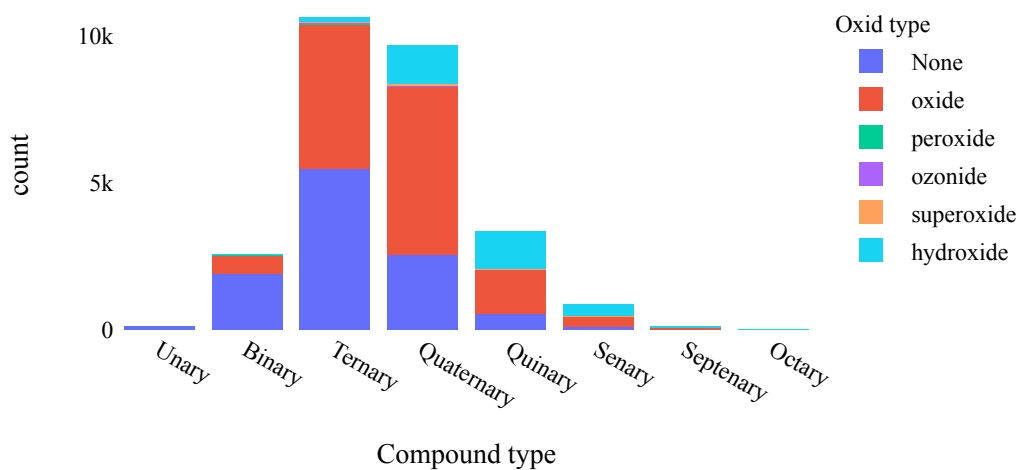


Figure 7.1: Distribution of oxid types as a function of number of elements in compounds in the data. The majority of the entries are found as oxides, while the second most frequent type is not an oxid.

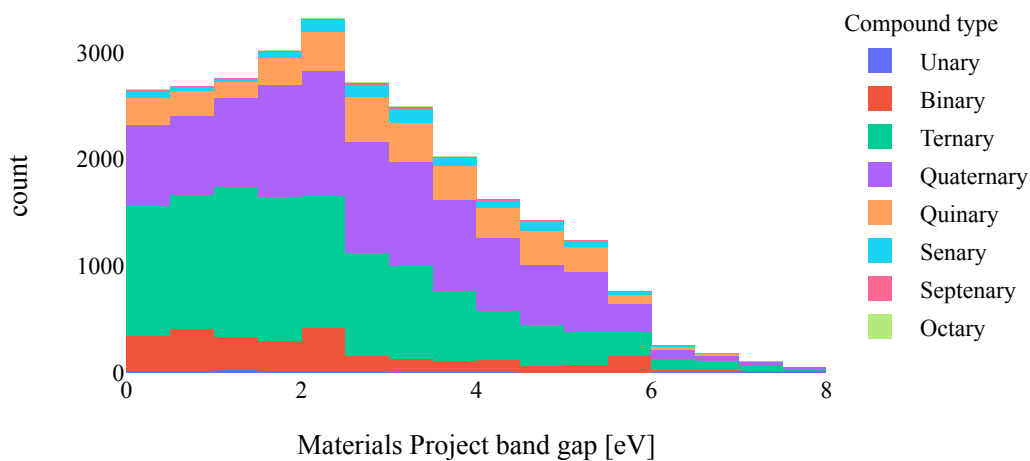


Figure 7.2: Distribution of band gaps as function of compound type in the data. The majority of compounds are ternary and quaternary, while the simpler compounds are few.

The second figure visualize the compound type as function of band gap, as calculated by Materials Project. Most of the materials existing in the data has a band gap lower than 2.3eV, where ternary compounds are most prominent. For larger values, we observe that quaternary compounds becomes dominant for larger values.

7.1 Comparing functionals for band gaps

Since the true size of a band gap can not be accurately determined by ab-initio calculations, we provide information regarding five different methods to obtain band gaps as visualized in figure 7.3. We have extracted experimental band gaps from Citrine Informatics that match the entries made by the initial MP query, involving entries that are associated with an ICSD structure that have a PBE-GGA calculated band gaps of minimum 0.1. All the band gaps to the left are found common with all databases through screening of correct structure, space group and ICSD-ID, while the figures to the right are only compared to the experimental database of Citrine Informatics. We found it helpful with the ICSD-tag to find similarities due to databases often have different norms and data-structures of descriptors, which proves challenging for comparison of stored calculations. If we were to exclude ICSD-tags, it would result in a much larger foundation to find similar entries, however, we found that the determination of similar entries would experience a large deviation when it comes to structures. By including an ICSD-tag, we reduce the basis of comparison but find more than 98% identical space group for entries in each database compared to Materials Project.

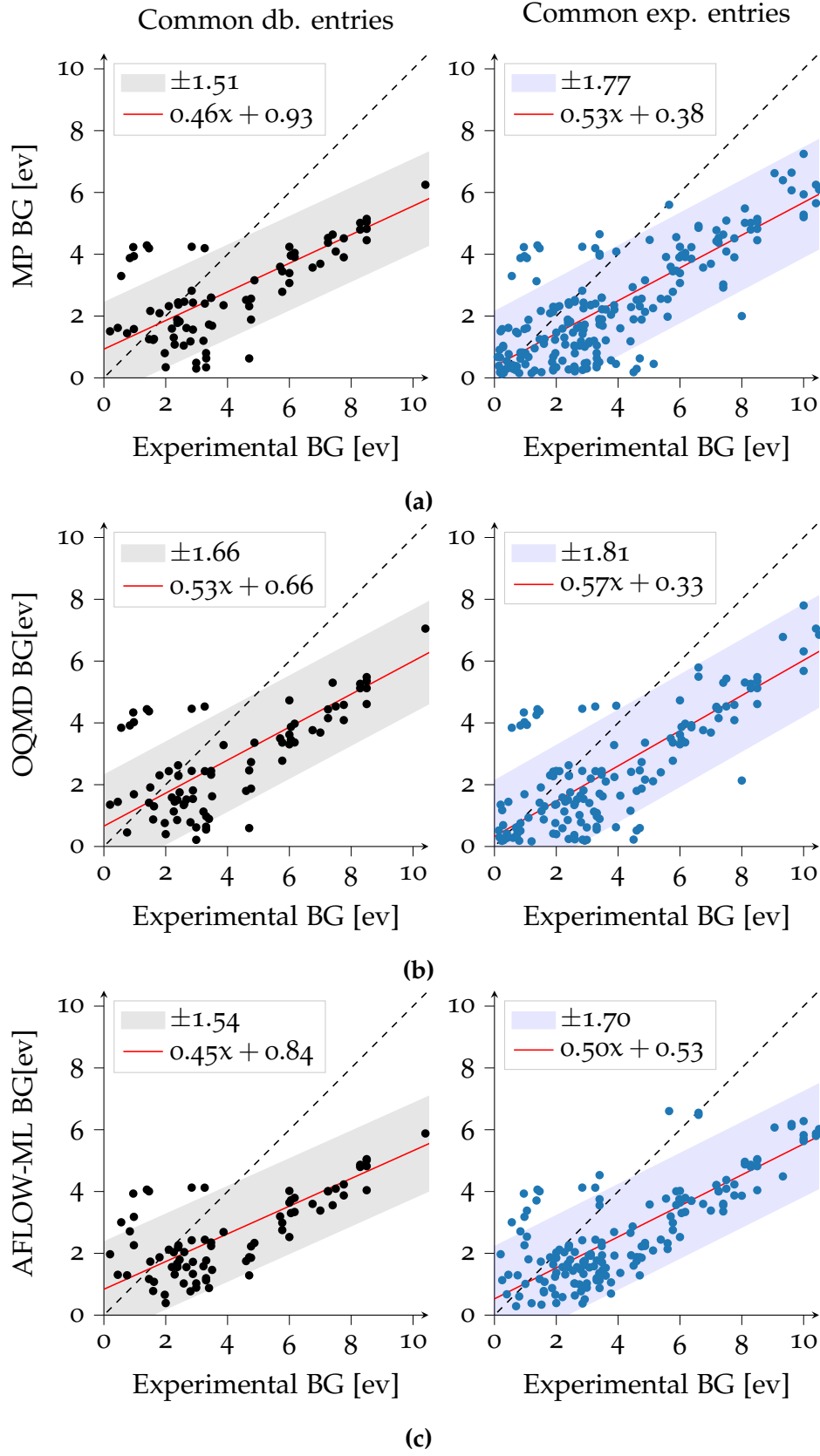
It was found that a very small portion of the data extracted from AFLOW was associated with an ICSD-tag, only 5 similar entries to the other databases, and therefore we have excluded the database from further consideration.

In the figures of 7.3, we observe each entry marked as black or blue dots. The dotted lines visualize the optimal ratio of estimated band gap to experimental values, while the red lines shows an linear least square fit to the data with the scrabbled area being the 95% confidence interval. The data that constitute the left figures are based on 82 similar entries, while the right figures constitute of more entries depending on the respective database. The data restriction was due to a small experimental database.

Initially, we wanted to include the right figures in the attempt of reducing the confidence interval with increasing the data points, but instead we find that the uncertainty of the confidence interval increase for all ab-initio calculations. This is due to the fact that the majority of the new entries are found for low band gap values, where the mismatch between experimental and calculated values are the largest. The discrepancy seems to be largest for values under 5eV, where entries are either calculated to have a very large band gap

where the experimental values report a very low band gap, which is also true the opposite case. One reason for this is that we have no information regarding the experiment where the band gap was determined. The information we from the experimental database is only considered the chemical formula of a compound, whereas the structure or ICSD-entry remains unknown. However, the same data of experimental values have been considered through other articles [101, 110].

Therefore, we find that the functional applied for Materials Project are found to underestimate the band gap with 30 – 60% while OQMD underestimates the band gap by 25 – 55%. AFLOW-ML also severely underestimates the band gap by 30 – 60%, but additionally have problems to accurately predict if a material is a metal or not. Many materials with both experimental and ab-initio calculations that showed a band gap of more than 1eV was predicted as metals by AFLOW-ML. JARVIS-DFT, on the other hand, was found to underestimate the band gap by 20 – 60% for the OptB88 and 0 – 30% for TB-mBJ functionals.



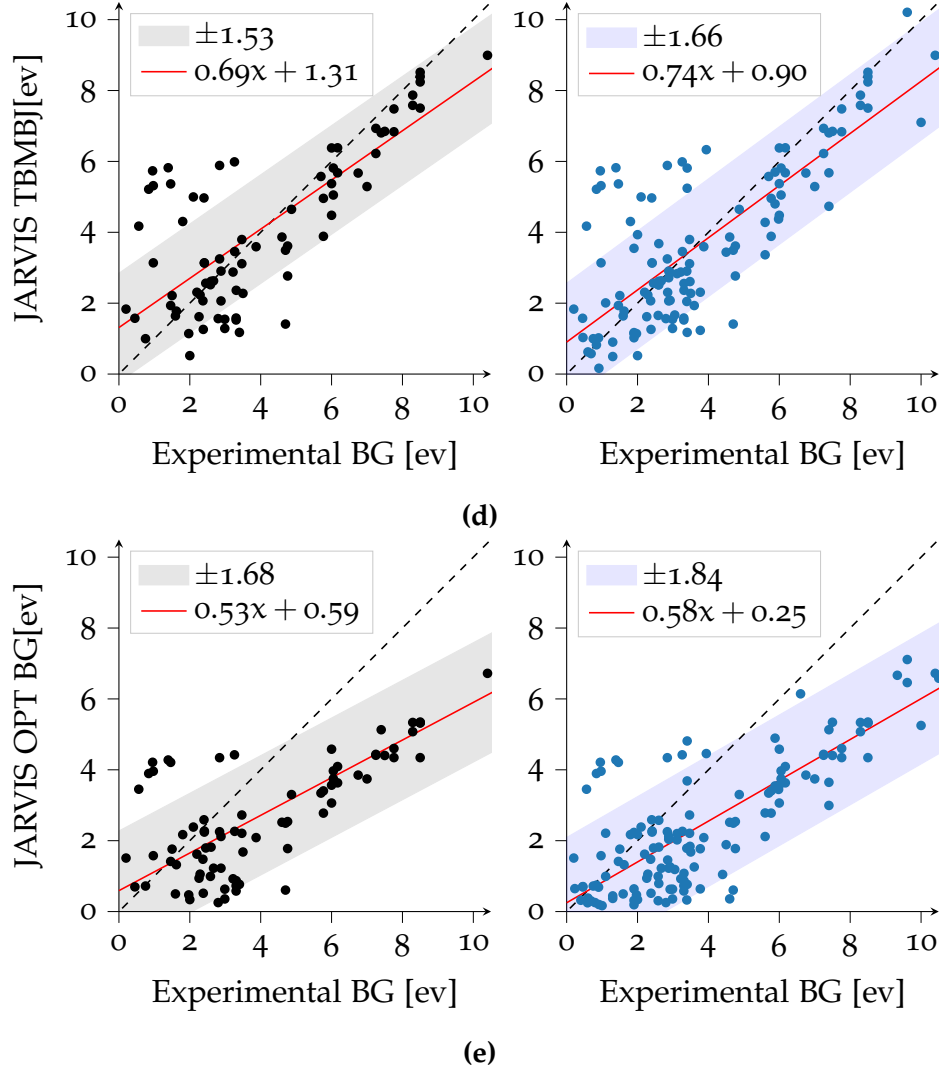


Figure 7.3: Comparison of reported experimental band gaps to those calculated by (a) Materials Project, (b) Open Quantum Materials Database, (c) AFLOW-ML, (d) JARVIS-DFT (TB-mBJ) and (e) JARVIS-DFT (OptB88). The figures to the left show reported band gaps that have been found to be common through all databases, while the figures to the right are only common with experimentally reported values from Citrine Informatics. All entries have been extracted in the period of january to march of 2021.

7.2 Technical details on ML classifiers

In the evaluation of the approaches, we apply a 5×5 stratified cross-validation when iterating through the hyperparameter combinations. For random forest, gradient boost and decision tree, we found that by adjusting the parameters responded to severe overfitting except for the default values defined by Scikit-learn. The only parameter that we found could potentially improve the evaluation metric f_1 was maximum number of depth for the trees grown, which we adjusted between 1 and 8. For logistic regression, we choose to adjust the regularisation strength with seven logarithmical adjusted values 10^{-3} to 10^5 , and use either 200 or 400 iterations to reach convergence.

When searching for the optimal number of principal components, we iterated over every odd number of principal components from 1 to the upper restricted number which defines an accumulated variance of 95% from the principal component analysis. Due to the large number of principal components, we end up fitting 25 folds for each of 1232 parameter combinations, totalling up to 30800 individual models, just for logistic regression for one approach. This serves as an additional motivator to keep the models simple, and accordingly shows how easy an initial complex step might evolve into an unfeasible amount of information. Therefore, we will not make an extensive analysis for every model, but emphasis important distinctions between the general models and provide background for principal choices made. However, it should be noted that a larger automated analysis is distributed through the MIT license at the Github repository *predicting-solid-state-qubit-candidates* [117].

7.2.1 The Ferrenti approach

We visualize the grid search for the optimal number of principal components in figure 7.4, where we present the mean training score, mean test score and mean f_1 score as a function of principal components used in the models. For each principal component, we visualize the optimal combination of hyperparameters based on the f_1 -score in the model. Common to all models is the improvement of scores up to around 50 principal components, where random forest and the decision tree slowly starts to overfit. From figure 7.4b, we see that the metrics gets heavily effected for larger values of principal components except for decision trees. We can observe that the f_1 -score is not varying as much as the other metrics due to an increasing number of positive predictions. This means that the accuracy of positive predictions are dominating the overall accuracy measurement, and we would expect a large amount of training data being predicted as positive candidates. Similar to decision trees are random forest, which also show sign of overfitting for larger values of principal components. However, as a result of several weak trees, it show smaller

signs of overfitting than the indications seen by the decision tree algorithm.

Gradient boost, on the other hand, experience minor changes for larger number of principal components, where the automated optimal component marked could easily be 50 principal components less without any remarks to the models metrics. We find that by using only a few principal components will make it reach almost 100% training accuracy, but does not show any sign on overfitting.

In 7.1, we find the precise measurements for each evaluation metric for the optimal number of principal components, which is visualized as dotted lines in figure 7.4. We find that the best performing model is logistic regression, but is dependent on a large amount of principal components. Random forest and gradient boost perform comparably, with and f1 score of 0.93 and 0.95, respectively.

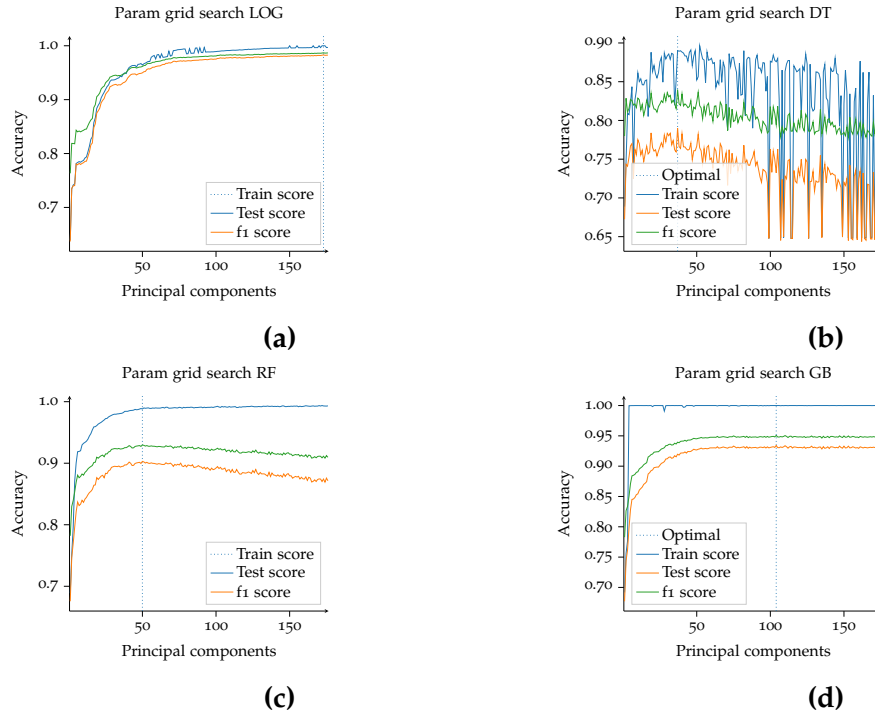


Figure 7.4: Four visualizations for each model of the grid-search with the resulting optimal number of principal components through a 5×5 stratified cross-validation. The upper limit of principal components is decided by the explained accumulated variance at 95%, while the optimal model is found by using the f1 score. The specific scores for the arbitrary number of principal components is found in the Appendix B.3.

Table 7.1: Table

Model	Optimal number of PC	Mean (std) test score	Mean (std) f1-score
LOG	173	0.98(0.008)	0.99(0.006)
DT	37	0.79(0.024)	0.84(0.020)
RF	50	0.90(0.020)	0.93(0.014)
GB	104	0.93(0.014)	0.95(0.010)

7.2.2 The augmented Ferrenti approach

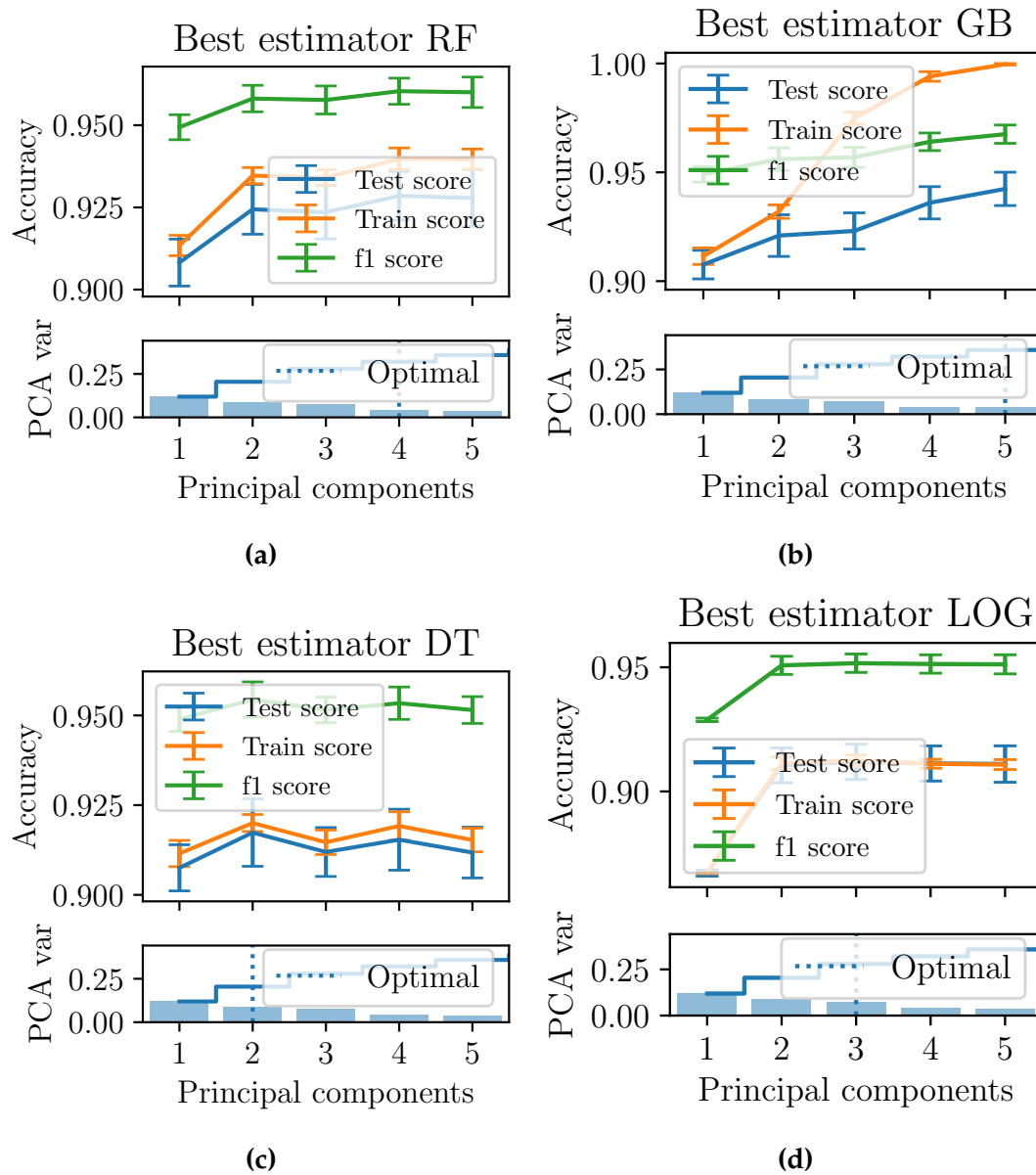


Figure 7.5: Four figures displaying hyperparameter search for the second approach. The best estimator is visualized for all hyperparameters as a function of principal components during a grid search with a 5x5 stratified cross validation. The lower plots visualizes the explained variance ratio, both accumulated and stepwise. The dotted lines marks the optimal hyperparameter-combination, while the error bars display the standard deviation.

7.2.3 The insightful approach

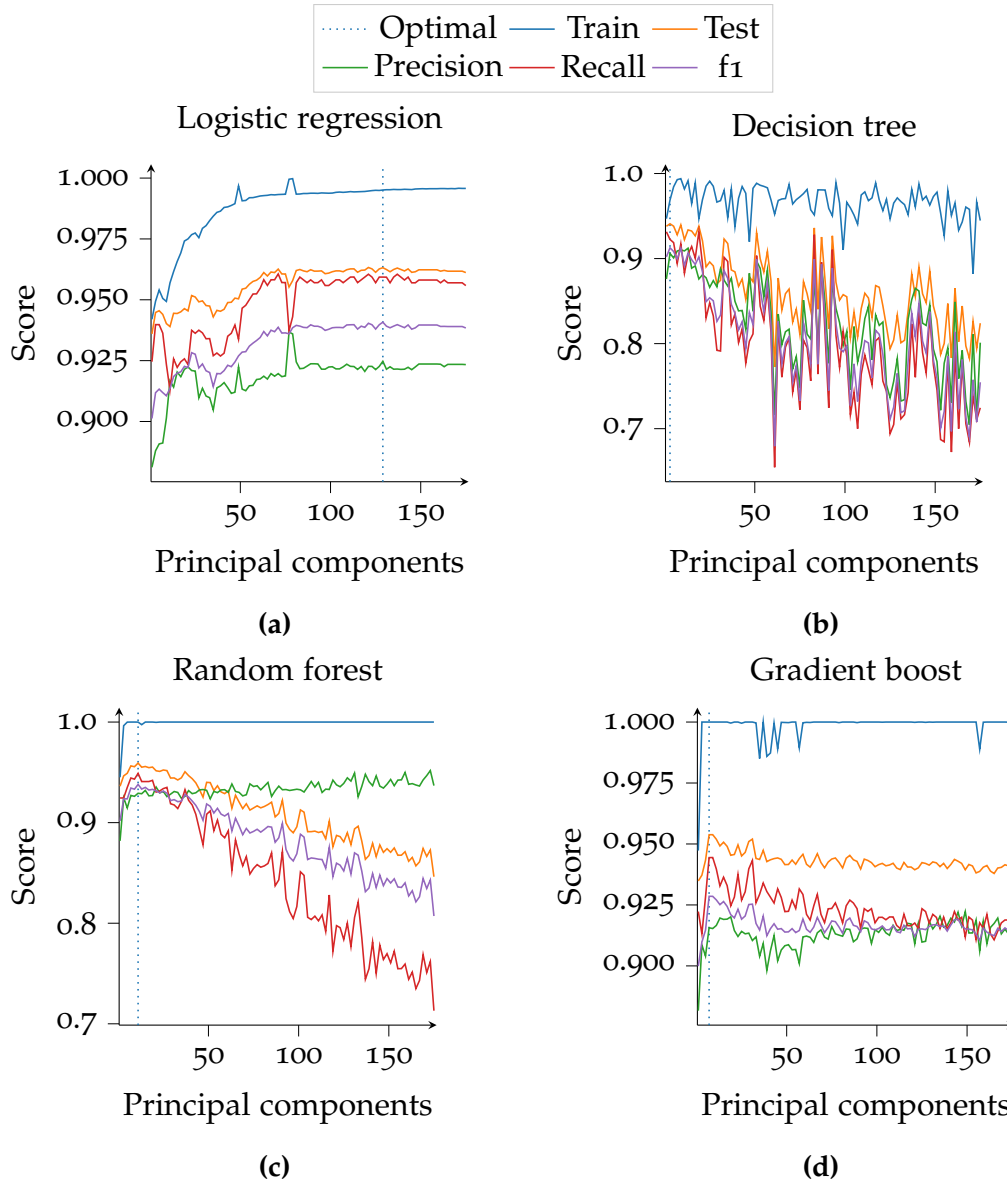


Figure 7.6: Four visualizations for each model of the grid-search with the resulting optimal number of principal components through a 5×5 stratified cross-validation. The upper limit of principal components is decided by the explained accumulated variance at 95%, while the optimal model is found using the f1 score.

Table 7.2: Table

Model	PC	Mean test	Mean recall	Mean precision	mean f1
LOG	129	0.96(0.018)	0.96(0.036)	0.93(0.041)	0.94(0.025)
DT	3	0.94(0.025)	0.92(0.050)	0.91(0.048)	0.91(0.032)
RF	11	0.96(0.019)	0.95(0.040)	0.93(0.039)	0.94(0.024)
GB	7	0.95(0.023)	0.94(0.047)	0.92(0.044)	0.93(0.030)

Part V

Appendices

Appendix A

Extracting information from ab-initio calculations

A.1 The Born-Oppenheimer approximation

The many-particle eigenfunction describes the wavefunction of all the electrons and nuclei and we denote it as Ψ_{κ}^{en} for electrons (e) and nuclei (n), respectively. The Born-Oppenheimer approximation assumes that nuclei, of substantially larger mass than electrons, can be treated as fixed point charges. According to this assumption, we can separate the eigenfunction into an electronic part and a nuclear part,

$$\Psi_{\kappa}^{en}(\mathbf{r}, \mathbf{R}) \approx \Psi_{\kappa}(\mathbf{r}, \mathbf{R}) \Theta_{\kappa}(\mathbf{R}), \quad (\text{A.1})$$

where the electronic part is dependent on the nuclei. This is in accordance with the assumption above, since electrons can respond instantaneously to a new position of the much slower nucleus, but this is not true for the opposite scenario. To our advantage, we already have knowledge of the terms in the many-particle Hamiltonian, and we can begin by separating the Hamiltonian into electronic and nuclear parts:

$$\hat{H}^{en} = \overbrace{T_e + U_{ee} + U_{en}}^{\hat{H}^e} + \overbrace{T_n + U_{nn}}^{\hat{H}^n}. \quad (\text{A.2})$$

Starting from the Schrödinger equation, we can formulate separate expressions for the electronic and the nuclear Schrödinger equations.

$$\hat{H}^{en}\Psi_{\kappa}^{en}(\mathbf{r}, \mathbf{R}) = E_{\kappa}^{en}\Psi_{\kappa}^{en}(\mathbf{r}, \mathbf{R}) \quad |\times \int \Psi^*(\mathbf{r}, \mathbf{R}) d\mathbf{r} \quad (\text{A.3})$$

$$\int \Psi_{\kappa}^*(\mathbf{r}, \mathbf{R})(\hat{H}^e + \hat{H}^n)\Psi_{\kappa}(\mathbf{r}, \mathbf{R})\Theta_{\kappa}(\mathbf{R}) d\mathbf{r} = E_{\kappa}^{en} \underbrace{\int \Psi_{\kappa}^*(\mathbf{r}, \mathbf{R})\Psi_{\kappa}(\mathbf{r}, \mathbf{R}) d\mathbf{r}}_1 \Theta_{\kappa}(\mathbf{R}). \quad (\text{A.4})$$

Since $\Theta_{\kappa}(\mathbf{R})$ is independent of the spatial coordinates to electrons, we get E_{κ} as the total energy of the electrons in the state κ .

$$E_{\kappa}(\mathbf{R})\Theta_{\kappa}(\mathbf{R}) + \int \Psi_{\kappa}^*(\mathbf{r}, \mathbf{R})\hat{H}^n\Psi_{\kappa}(\mathbf{r}, \mathbf{R})\Theta_{\kappa}(\mathbf{R}) d\mathbf{r} = E_{\kappa}^{en}\Theta_{\kappa}(\mathbf{R}). \quad (\text{A.5})$$

Now, the final integration term can be simplified by using the product rule, which results in

$$\left(T_n + T_n' + T_n'' + U_{nn} + E_{\kappa}(\mathbf{R})\right)\Theta_{\kappa}(\mathbf{R}) = E_{\kappa}^{en}\Theta_{\kappa}(\mathbf{R}). \quad (\text{A.6})$$

If we neglect T_n' and T_n'' to lower the computational efforts, we obtain the Born-Oppenheimer approximation with the electronic eigenfunction as

$$(T_e + U_{ee} + U_{en})\Psi_{\kappa}(\mathbf{r}, \mathbf{R}) = E_{\kappa}(\mathbf{R})\Psi_{\kappa}(\mathbf{r}, \mathbf{R}) \quad (\text{A.7})$$

and the nuclear eigenfunction as

$$(T_n + U_{nn} + E_{\kappa}(\mathbf{R}))\Theta_{\kappa}(\mathbf{R}) = E_{\kappa}^{en}(\mathbf{R})\Theta_{\kappa}(\mathbf{R}). \quad (\text{A.8})$$

How are they coupled, you might ask? The total energy in the electronic equation is a potential in the nuclear equation.

A.2 The variational principle

So far, we have tried to make the time-independent Schrödinger equation easier with the use of an *ansatz*, but we do not necessarily have an adequate guess for the eigenfunctions and the ansatz can only give a rough estimate in most scenarios. Another approach, namely the *variational principle*, states that the energy of any trial wavefunction is always an upper bound to the exact ground state energy by definition E_0 .

$$E_0 = \langle \psi_0 | H | \psi_0 \rangle \leq \langle \psi | H | \psi \rangle = E \quad (\text{A.9})$$

The eigenfunctions of H form a complete set, which means any normalized Ψ can be expressed in terms of the eigenstates

$$\Psi = \sum_n c_n \psi_n, \quad \text{where} \quad H\psi_n = E_n \psi_n \quad (\text{A.10})$$

for all $n = 1, 2, \dots$. The expectation value for the energy can be calculated as

$$\begin{aligned} \langle \Psi | H | \Psi \rangle &= \left\langle \sum_n c_n \psi_n \left| H \right| \sum_{n'} c_{n'} \psi_{n'} \right\rangle \\ &= \sum_n \sum_{n'} c_n^* c_{n'} \langle \psi_n | H | \psi_{n'} \rangle \\ &= \sum_n \sum_{n'} c_n^* c_{n'} E_n \langle \psi_n | \psi_{n'} \rangle \end{aligned}$$

Here we assume that the eigenfunctions have been orthonormalized and we can utilize $\langle \psi_m | \psi_n \rangle = \delta_{mn}$, resulting in

$$\sum_n c_n^* c_n E_n = \sum_n |c_n|^2 E_n.$$

We have already stated that Ψ is normalized, thus $\sum_n |c_n|^2 = 1$, and the expectation value conveniently is bound to follow equation A.9. The quest to understand the variational principle can be summarized in a sentence - it is possible to tweak the wavefunction parameters to minimize the energy, or summed up in a mathematical phrase,

$$E_0 = \min_{\Psi \rightarrow \Psi_0} \langle \Psi | H | \Psi \rangle. \quad (\text{A.11})$$

A.3 The Hohenberg-Kohn theorems

A.3.1 The Hohenberg-Kohn theorem 1

PROOF. Assume that two external potentials $V_{\text{ext}}^{(1)}$ and $V_{\text{ext}}^{(2)}$, that differ by more than a constant, have the same ground state density $n_0(\mathbf{r})$. The two different potentials correspond to distinct Hamiltonians $\hat{H}_{\text{ext}}^{(1)}$ and $\hat{H}_{\text{ext}}^{(2)}$, which again give rise to distinct wavefunctions $\Psi_{\text{ext}}^{(1)}$ and $\Psi_{\text{ext}}^{(2)}$. Utilizing the variational principle, we find that no wavefunction can give an energy that is less than the energy of $\Psi_{\text{ext}}^{(1)}$ for $\hat{H}_{\text{ext}}^{(1)}$, that is

$$E^{(1)} = \langle \Psi^{(1)} | \hat{H}^{(1)} | \Psi^{(1)} \rangle < \langle \Psi^{(2)} | \hat{H}^{(1)} | \Psi^{(2)} \rangle \quad (\text{A.12})$$

and

$$E^{(2)} = \langle \Psi^{(2)} | \hat{H}^{(2)} | \Psi^{(2)} \rangle < \langle \Psi^{(1)} | \hat{H}^{(2)} | \Psi^{(1)} \rangle. \quad (\text{A.13})$$

Assuming that the ground state is not degenerate, the inequality strictly holds. Since we have identical ground state densities for the two Hamiltonian's, we can rewrite the expectation value for equation A.12 as

$$\begin{aligned} E^{(1)} &= \langle \Psi^{(1)} | \hat{H}^{(1)} | \Psi^{(1)} \rangle \\ &= \langle \Psi^{(1)} | T + U_{ee} + U_{\text{ext}}^{(1)} | \Psi^{(1)} \rangle \\ &= \langle \Psi^{(1)} | T + U_{ee} | \Psi^{(1)} \rangle + \int \Psi^{*(1)}(\mathbf{r}) V_{\text{ext}}^{(1)}(\mathbf{r}) \Psi^{(1)}(\mathbf{r}) d\mathbf{r} \\ &= \langle \Psi^{(1)} | T + U_{ee} | \Psi^{(1)} \rangle + \int V_{\text{ext}}^{(1)}(\mathbf{r}) n(\mathbf{r}) d\mathbf{r} \\ &< \langle \Psi^{(2)} | \hat{H}^{(1)} | \Psi^{(2)} \rangle \\ &= \langle \Psi^{(2)} | T + U_{ee} + U_{\text{ext}}^{(1)} + \overbrace{U_{\text{ext}}^{(2)} - U_{\text{ext}}^{(2)}}^0 | \Psi^{(2)} \rangle \\ &= \langle \Psi^{(2)} | T + U_{ee} + U_{\text{ext}}^{(2)} | \Psi^{(1)} \rangle + \int (V_{\text{ext}}^{(1)} - V_{\text{ext}}^{(2)}) n(\mathbf{r}) d\mathbf{r} \\ &= E^{(2)} + \int (V_{\text{ext}}^{(1)} - V_{\text{ext}}^{(2)}) n(\mathbf{r}) d\mathbf{r}. \end{aligned}$$

Thus,

$$E^{(1)} = E^{(2)} + \int (V_{\text{ext}}^{(1)} - V_{\text{ext}}^{(2)}) n(\mathbf{r}) d\mathbf{r} \quad (\text{A.14})$$

A similar procedure can be performed for $E^{(2)}$ in equation A.13, resulting in

$$E^{(2)} = E^{(1)} + \int (V_{\text{ext}}^{(2)} - V_{\text{ext}}^{(1)}) n(\mathbf{r}) d\mathbf{r}. \quad (\text{A.15})$$

If we add these two equations together, we get

$$\begin{aligned} E^{(1)} + E^{(2)} &< E^{(2)} + E^{(1)} + \int (V_{\text{ext}}^{(1)} - V_{\text{ext}}^{(2)}) n(\mathbf{r}) d\mathbf{r} \\ &\quad + \int (V_{\text{ext}}^{(2)} - V_{\text{ext}}^{(1)}) n(\mathbf{r}) d\mathbf{r} \\ E^{(1)} + E^{(2)} &< E^{(2)} + E^{(1)}, \end{aligned} \quad (\text{A.16})$$

which is a contradiction. Thus, the two external potentials cannot have the same ground-state density, and $V_{\text{ext}}(\mathbf{r})$ is determined uniquely (except for a constant) by $n(\mathbf{r})$. \square

A.3.2 The Hohenberg-Kohn theorem 2

PROOF. Since the external potential is uniquely determined by the density and since the potential in turn uniquely determines the ground state wavefunction (except in degenerate situations), all the other observables of the system are uniquely determined. Then the energy can be expressed as a functional of the density.

$$E[n] = \underbrace{T[n] + U_{ee}[n]}_{F[n]} + \underbrace{\int V_{en}n(r)dr}_{U_{en}[n]} \quad (\text{A.17})$$

where $F[n]$ is a universal functional because the treatment of the kinetic and internal potential energies are the same for all systems, however, it is most commonly known as the Hohenberg-Kohn functional.

In the ground state, the energy is defined by the unique ground-state density $n_0(r)$,

$$E_0 = E[n_0] = \langle \Psi_0 | H | \Psi_0 \rangle. \quad (\text{A.18})$$

From the variational principle, a different density $n(r)$ will give a higher energy

$$E_0 = E[n_0] = \langle \Psi_0 | H | \Psi_0 \rangle < \langle \Psi | H | \Psi \rangle = E[n] \quad (\text{A.19})$$

Thus, the total energy is minimized for n_0 , and so has to be the ground-state energy. \square

A.4 Self-consistent field methods

So, the remaining question is, how do we solve the Kohn-Sham equation? First, we would need to define the Hartree potential, which can be found if we know the electron density. The electron density can be found from the single-electron wave-functions, however, these can only be found from solving the Kohn-Sham equation. This *circle of life* has to start somewhere, but where? The process can be defined as an iterative method, a *computational scheme*, as visualized in figure A.1.

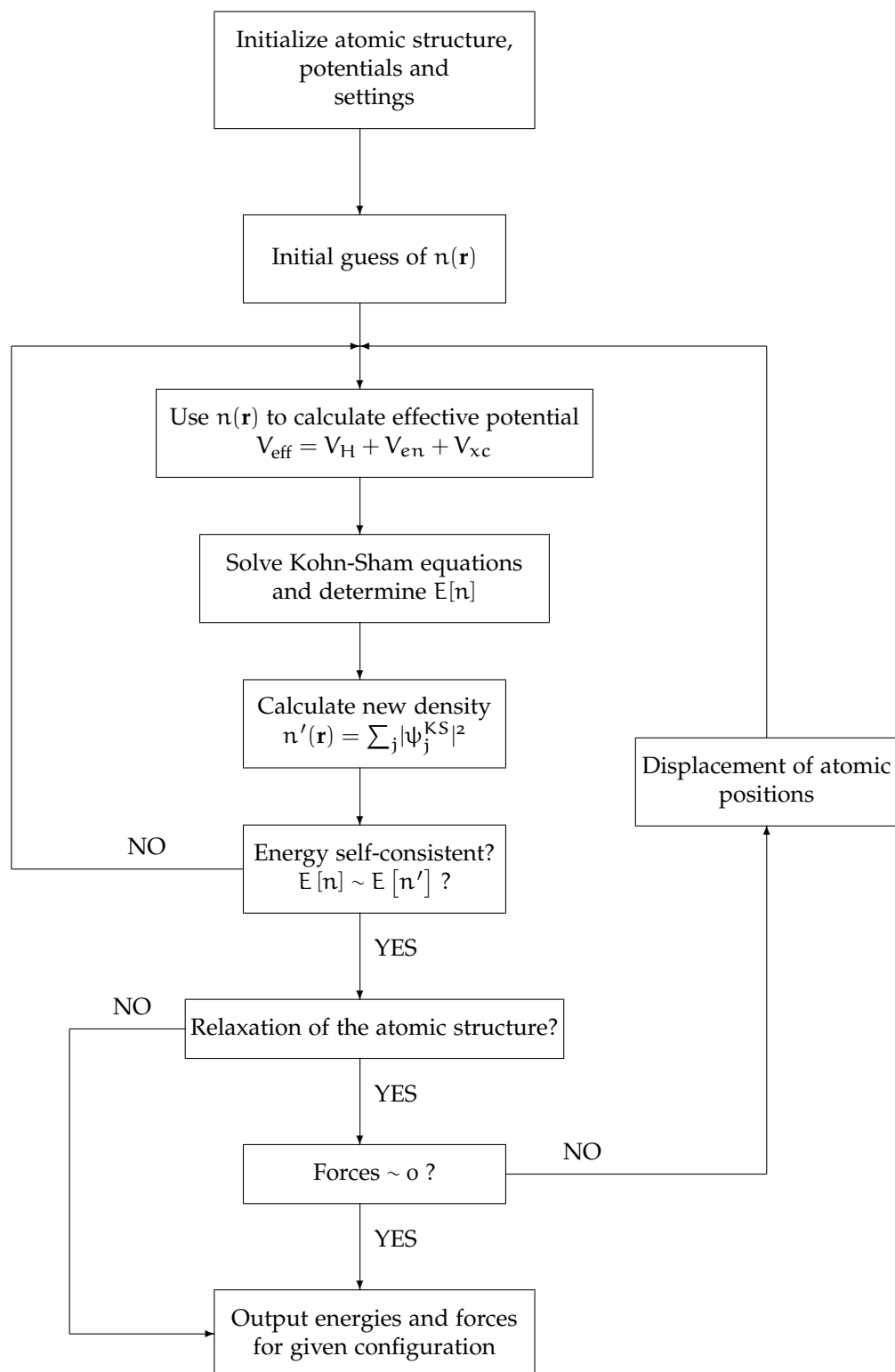


Figure A.1: A flow chart of the self-consistent field method for DFT.

Appendix B

Featurization

B.1 Table of featurizers

Table B.1: This thesis' chosen 39 featurizers from matminer. Descriptions are either found from Ref. [101] or from the project's Github page.

Features	Description	Original reference
Composition features		
AtomicOrbitals	Highest occupied molecular orbital (HOMO) and lowest unoccupied molecular orbital (LUMO).	[128]
AtomicPacking-Efficiency	Packing efficiency.	[129]
BandCenter	Estimation of absolute position of band center using geometric mean of electronegativity.	[130]
ElementFraction	Fraction of each element in a composition.	-
ElementProperty	Statistics of various element properties.	[100, 131, 132]
IonProperty	Maximum and average ionic character.	[131]
Continued on next page		

Table B.1 – continued from previous page

Features	Description	Original reference
Miedema	Formation enthalpies of intermetallic compounds, solid solutions, and amorphous phases using semi-empirical Miedema model.	[133]
Stoichiometry	L^p norm-based stoichiometric attributes.	[131]
TMetalFraction	Fraction of magnetic transition metals.	[132]
ValenceOrbital	Valence orbital attributes such as the mean number of electrons in each shell.	[131]
YangSolid-Solution	Mixing thermochemistry and size mismatch terms.	[134]
Oxid composition features		
Electronegativity-Diff	Statistics on electronegativity difference between anions and cations.	[132]
OxidationStates	Statistics of oxidation states.	[132]
Structure features		
DensityFeatures	Calculate density, volume per atom and packing fraction.	-
GlobalSymmetry-Features	Determines spacegroup number, crystal system (1-7) and inversion symmetry.	-
RadialDistribution-Function	Calculates the radial distribution function of a crystal system.	-
CoulombMatrix	Generate the Coulomb matrix, which is a representation of the nuclear coulombic interaction of the input structure.	[135]
PartialRadial-Distribution-Function	Compute the partial radial distribution function of a crystal structure	[136]

Continued on next page

Table B.1 – continued from previous page

Features	Description	Original reference
SineCoulomb-Matrix	Computes a variant of the coulomb matrix developed for periodic crystals.	[137]
EwaldEnergy	Computes the energy from Coulombic interactions based on charge states of each site.	[138]
BondFractions	Compute the fraction of each bond in a structure, based on nearest neighbours.	[139]
Structural-Heterogeneity	Calculates the variance in bond lengths and atomic volumes in a structure.	[140]
MaximumPacking-Efficiency	Calculates the maximum packing efficiency of a structure.	[140]
ChemicalOrdering	ordering of species differs from random in a structure.	[140]
XRDPowder-Pattern	1D array representing normalized powder diffraction of a structure as calculated by pymatgen.	[100]
Site features		
AGNI-Fingerprints	Calculates the product integral of RDF and Gaussian window function	[141]
AverageBond-Angle	Determines the average bond angle of a specific site with its nearest neighbors using pymatgens implementation.	[142]
AverageBond-Length	Determines the average bond length between one specific site and all its nearest neighbors using pymatgens implementation.	[142]
BondOrientational-Paramater	Calculates the averages of spherical harmonics of local neighbors	[143, 144]
Continued on next page		

Table B.1 – continued from previous page

Features	Description	Original reference
ChemEnvSite Fingerprint	Calculates the resemblance of given sites to ideal environment using pymatgens ChemEnv package.	[145, 146]
Coordination- Number	The number of first nearest neighbors of a site	[146]
CrystalNN- Fingerprint	A local order parameter fingerprint for periodic crystals.	-
GaussianSymm- Func	Calculates the gaussian radial and angular symmetry functions originally suggested for fitting machine learning potentials.	[147, 148]
GeneralizedRadial- Distribution- Function	Computes the general radial distribution function for a site	[143]
LocalProperty- Difference	Computes the difference in elemental properties between a site and its neighboring sites.	[140, 142]
OPSite- Fingerprint	Computes the local structure order parameters from a site's neighbor environment.	[146]
Voronoi- Fingerprint	Calculates the Voronoi tessellation-based features around a target site.	[149, 150]
Density of state features		
DOSFeaturizer	Computes top contributors to the density of states at the valence and conduction band edges. Thus includes chemical specie, orbital character, and orbital location information.	[151]
Band structure features		
Continued on next page		

Table B.1 – continued from previous page

Features	Description	Original reference
BandFeaturizer	Converts a complex electronic band structure into quantities such as band gap and the norm of k point coordinates at which the conduction band minimum and valence band maximum occur.	-

B.2 Erroneous entries

MPID	Full formula	Reference
mp-555563	$\text{PH}_6\text{C}_2\text{S}_2\text{NCl}_2\text{O}_4$	[152]
mp-583476	$\text{Nb}_7\text{S}_2\text{I}_{19}$	[153]
mp-600205	$\text{H}_{10}\text{C}_5\text{SeS}_2\text{N}_3\text{Cl}$	-
mp-600217	$\text{H}_{80}\text{C}_{40}\text{Se}_8\text{S}_{16}\text{Br}_8\text{N}_{24}$	-
mp-1195290	$\text{Ga}_3\text{Si}_5\text{P}_{10}\text{H}_{36}\text{C}_{12}\text{N}_4\text{Cl}_{11}$	-
mp-1196358	$\text{P}_4\text{H}_{120}\text{Pt}_8\text{C}_{40}\text{I}_8\text{N}_4\text{Cl}_8$	-
mp-1196439	$\text{Sn}_8\text{P}_4\text{H}_{128}\text{C}_{44}\text{N}_{12}\text{Cl}_8\text{O}_4$	-
mp-1198652	$\text{Te}_4\text{H}_{72}\text{C}_{36}\text{S}_{24}\text{N}_{12}\text{Cl}_4$	-
mp-1198926	$\text{Re}_8\text{H}_{96}\text{C}_{24}\text{S}_{24}\text{N}_{48}\text{Cl}_{48}$	-
mp-1199490	$\text{Mn}_4\text{H}_{64}\text{C}_{16}\text{S}_{16}\text{N}_{32}\text{Cl}_8$	-
mp-1199686	$\text{Mo}_4\text{P}_{16}\text{H}_{152}\text{C}_{52}\text{N}_{16}\text{Cl}_{16}$	-
mp-1203403	$\text{C}_{121}\text{S}_2\text{Cl}_{20}$	-
mp-1204279	$\text{Si}_{16}\text{Te}_8\text{H}_{176}\text{Pd}_8\text{C}_{64}\text{Cl}_{16}$	-
mp-1204629	$\text{P}_{16}\text{H}_{216}\text{C}_{80}\text{N}_{32}\text{Cl}_8$	-

Table B.2: A table of manually identified entries from Materials Project that experience issues concerning Matminer’s featurization tools. These were excluded from the dataset.

B.3 Optimalization

Table B.3: Table of scores versus principal components used in the model for the Ferrenti approach.

Model	PC	Mean (std) test score	Mean (std) f1-score
LOG	1	0.64(0.010)	0.76(0.007)
	2	0.74(0.020)	0.82(0.013)
	5	0.78(0.018)	0.84(0.013)
	10	0.79(0.019)	0.85(0.013)
	50	0.95(0.014)	0.96(0.010)
	176	0.98(0.010)	0.99(0.008)
Optimal	173	0.98(0.008)	0.99(0.006)
DT	1	0.67(0.015)	0.78(0.011)
	2	0.74(0.017)	0.83(0.010)
	5	0.77(0.024)	0.83(0.019)
	10	0.76(0.028)	0.82(0.024)
	50	0.77(0.029)	0.82(0.021)
	176	0.72(0.030)	0.79(0.023)
Optimal	37	0.79(0.024)	0.84(0.020)
RF	1	0.68(0.015)	0.78(0.010)
	2	0.75(0.019)	0.83(0.011)
	5	0.82(0.023)	0.87(0.016)
	10	0.84(0.025)	0.88(0.017)
	50	0.90(0.020)	0.93(0.014)
	176	0.87(0.016)	0.91(0.010)
Optimal	50	0.90(0.020)	0.93(0.014)
GB	1	0.68(0.015)	0.78(0.010)
	2	0.74(0.019)	0.83(0.011)
	5	0.83(0.020)	0.87(0.015)
	10	0.85(0.020)	0.89(0.015)
	50	0.93(0.017)	0.95(0.012)
	176	0.93(0.017)	0.95(0.012)
Optimal	104	0.93(0.014)	0.95(0.010)

Bibliography

1. Griffiths, D. *Introduction to quantum mechanics* ISBN: 9781107179868 (Cambridge University Press, Cambridge, 2017).
2. Turing, A. M. On computable numbers, with an application to the Entscheidungsproblem. *Proceedings of the London mathematical society* **2**, 230–265 (1937).
3. Moore, G. Cramming More Components Onto Integrated Circuits. *Proceedings of the IEEE* **86**, 82–85 (Jan. 1965).
4. Pavičić, M. *Quantum computation and quantum communication : theory and experiments* ISBN: 9786610743704 (Springer, New York, 2006).
5. Gwennap, L. Apple's 5 Nanometer Chip Is Another Signpost That Moore's Law Is Running Out. *Forbes*. <<https://www.forbes.com/sites/linleygwennap/2020/10/12/apple-moores-law-is-running-out/>> (Oct. 12, 2020).
6. Weber, J. R. *et al.* Quantum computing with defects. *Proceedings of the National Academy of Sciences* **107**, 8513–8518 (Apr. 2010).
7. DiVincenzo, D. P. The Physical Implementation of Quantum Computation. *Fortschritte der Physik* **48**, 771–783 (Sept. 2000).
8. Ladd, T. D. *et al.* Quantum computers. *Nature* **464**, 45–53 (Mar. 2010).
9. Mizel, A., Lidar, D. A. & Mitchell, M. Simple Proof of Equivalence between Adiabatic Quantum Computation and the Circuit Model. *Physical Review Letters* **99**. doi:[10.1103/physrevlett.99.070502](https://doi.org/10.1103/physrevlett.99.070502) (Aug. 2007).
10. Grover, L. K. A framework for fast quantum mechanical algorithms. arXiv: [quant-ph/9711043v2](https://arxiv.org/abs/quant-ph/9711043v2) [[quant-ph](https://arxiv.org/abs/quant-ph)] (Nov. 20, 1997).
11. Shor, P. *Algorithms for quantum computation: discrete logarithms and factoring* in *Proceedings 35th Annual Symposium on Foundations of Computer Science* (IEEE Comput. Soc. Press, 1994). doi:[10.1109/sfcs.1994.365700](https://doi.org/10.1109/sfcs.1994.365700).
12. Martinis, J. M. *et al.* Quantum supremacy using a programmable superconducting processor en. 2019. doi:[10.5061/DRYAD.K6T1RJ8](https://doi.org/10.5061/DRYAD.K6T1RJ8).
13. Georgescu, I. The DiVincenzo criteria 20 years on. *Nature Reviews Physics* **2**, 666–666 (Nov. 2020).

14. Griffiths, R. B. Nature and location of quantum information. *Physical Review A* **66**. doi:[10.1103/physreva.66.012311](https://doi.org/10.1103/physreva.66.012311) (July 2002).
15. Gisin, N., Ribordy, G., Tittel, W. & Zbinden, H. Quantum cryptography. *Reviews of Modern Physics* **74**, 145–195 (Mar. 2002).
16. Gisin, N. & Thew, R. Quantum communication. *Nature Photonics* **1**, 165–171 (Mar. 2007).
17. Acín, A. *et al.* The quantum technologies roadmap: a European community view. *New Journal of Physics* **20**, 080201 (Aug. 2018).
18. Boaron, A. *et al.* Secure Quantum Key Distribution over 421 km of Optical Fiber. *Physical Review Letters* **121**. doi:[10.1103/physrevlett.121.190502](https://doi.org/10.1103/physrevlett.121.190502) (Nov. 2018).
19. Degen, C. L., Reinhard, F. & Cappellaro, P. Quantum sensing. *Reviews of Modern Physics* **89**. doi:[10.1103/revmodphys.89.035002](https://doi.org/10.1103/revmodphys.89.035002) (July 2017).
20. Kristian Fossheim, A. S. *Superconductivity: Physics and Applications* 442 pp. ISBN: 0470844523. <https://www.ebook.de/de/product/3608091/kristian_fossheim_asle_sudboe_superconductivity_physics_and_applications.html> (WILEY, 2004).
21. Ben Streetman, S. B. *Solid State Electronic Devices, Global Edition* 632 pp. ISBN: 1292060557. <https://www.ebook.de/de/product/30394493/ben_streetman_sanjay_banerjee_solid_state_electronic_devices_global_edition.html> (Pearson Education Limited, 2015).
22. Renganathan, G., Tanneru, N. & Madurai, S. L. in *Fundamental Biomaterials: Metals* 211–241 (Elsevier, 2018). doi:[10.1016/b978-0-08-102205-4.00010-6](https://doi.org/10.1016/b978-0-08-102205-4.00010-6).
23. Lufaso, M. W. & Woodward, P. M. Prediction of the crystal structures of perovskites using the software program SPuDS. *Acta Crystallographica Section B Structural Science* **57**, 725–738 (Nov. 2001).
24. Bednorz, J. G. & Müller, K. A. Perovskite-type oxides—The new approach to high-Tc superconductivity. *Reviews of Modern Physics* **60**, 585–600 (July 1988).
25. Boivin, J. C. & Mairesse, G. Recent Material Developments in Fast Oxide Ion Conductors. *Chemistry of Materials* **10**, 2870–2888 (Oct. 1998).
26. Cheong, S.-W. & Mostovoy, M. Multiferroics: a magnetic twist for ferroelectricity. *Nature Materials* **6**, 13–20 (Jan. 2007).
27. Ibn-Mohammed, T. *et al.* Perovskite solar cells: An integrated hybrid lifecycle assessment and review in comparison with other photovoltaic technologies. *Renewable and Sustainable Energy Reviews* **80**, 1321–1344 (Dec. 2017).

28. Chen, P.-Y. *et al.* Environmentally responsible fabrication of efficient perovskite solar cells from recycled car batteries. *Energy Environ. Sci.* **7**, 3659–3665 (2014).
29. Pauli, W. Über den Zusammenhang des Abschlusses der Elektronengruppen im Atom mit der Komplexstruktur der Spektren. *Zeitschrift für Physik* **31**, 765–783 (Feb. 1925).
30. Martienssen, W. *Springer handbook of condensed matter and materials data* ISBN: 9786610625949 (Springer, Heidelberg New York, 2005).
31. Pelant, I. *Luminescence spectroscopy of semiconductors* ISBN: 0191738549 (Oxford University Press, Oxford, 2012).
32. Kun Huang, A. R. Theory of light absorption and non-radiative transitions in F-centres. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences* **204**, 406–423 (Dec. 1950).
33. Gordon, L. *et al.* Quantum computing with defects. *MRS Bulletin* **38**, 802–807 (Oct. 2013).
34. Bernien, H. *et al.* Heralded entanglement between solid-state qubits separated by three metres. *Nature* **497**, 86–90 (Apr. 2013).
35. Taylor, J. M. *et al.* High-sensitivity diamond magnetometer with nanoscale resolution. *Nature Physics* **4**, 810–816 (Sept. 2008).
36. Barclay, P. E., Fu, K.-M. C., Santori, C., Faraon, A. & Beausoleil, R. G. Hybrid Nanocavity Resonant Enhancement of Color Center Emission in Diamond. *Physical Review X* **1**. doi:[10.1103/physrevx.1.011007](https://doi.org/10.1103/physrevx.1.011007) (Sept. 2011).
37. Neudeck, P. G. Progress in silicon carbide semiconductor electronics technology. *Journal of Electronic Materials* **24**, 283–288 (Apr. 1995).
38. Silveira, E., Freitas, J. A., Glembocki, O. J., Slack, G. A. & Schowalter, L. J. Excitonic structure of bulk AlN from optical reflectivity and cathodoluminescence measurements. *Physical Review B* **71**. doi:[10.1103/physrevb.71.041201](https://doi.org/10.1103/physrevb.71.041201) (Jan. 2005).
39. Lawaetz, P. Valence-Band Parameters in Cubic Semiconductors. *Physical Review B* **4**, 3460–3467 (Nov. 1971).
40. Beckers, L. *et al.* Structural and optical characterization of epitaxial waveguiding BaTiO₃ thin films on MgO. *Journal of Applied Physics* **83**, 3305–3310 (Mar. 1998).
41. Kumbhojkar, N., Nikesh, V. V., Kshirsagar, A. & Mahamuni, S. Photo-physical properties of ZnS nanoclusters. *Journal of Applied Physics* **88**, 6260–6264 (Dec. 2000).

42. Bassett, L. C., Alkauskas, A., Exarhos, A. L. & Fu, K.-M. C. Quantum defects by design. *Nanophotonics* **8**, 1867–1888 (Oct. 2019).
43. James, W. J. Theory of defects in solids by A. M. Stoneham. *Acta Crystallographica Section A* **32**, 527–527 (May 1976).
44. Togan, E. *et al.* Quantum entanglement between an optical photon and a solid-state spin qubit. *Nature* **466**, 730–734 (Aug. 2010).
45. Bathen, M. E. *Point defects in silicon carbide for quantum technologies: Identification, tuning and control* PhD thesis (The Faculty of Mathematics and Natural Sciences, University of Oslo).
46. Son, N. T. *et al.* Developing silicon carbide for quantum spintronics. *Applied Physics Letters* **116**, 190501 (May 2020).
47. Falk, A. L. *et al.* Polytype control of spin qubits in silicon carbide. *Nature Communications* **4**. doi:[10.1038/ncomms2854](https://doi.org/10.1038/ncomms2854) (May 2013).
48. Widmann, M. *et al.* Coherent control of single spins in silicon carbide at room temperature. *Nature Materials* **14**, 164–168 (Dec. 2014).
49. Zhang, G., Cheng, Y., Chou, J.-P. & Gali, A. Material platforms for defect qubits and single-photon emitters. *Applied Physics Reviews* **7**, 031308 (Sept. 2020).
50. Redjem, W. *et al.* Single artificial atoms in silicon emitting at telecom wavelengths. *Nature Electronics* **3**, 738–743 (Nov. 2020).
51. Wang, J. *et al.* Gallium arsenide (GaAs) quantum photonic waveguide circuits. *Optics Communications* **327**, 49–55 (Sept. 2014).
52. Berhane, A. M. *et al.* Photophysics of GaN single-photon emitters in the visible spectral range. *Physical Review B* **97**. doi:[10.1103/physrevb.97.165202](https://doi.org/10.1103/physrevb.97.165202) (Apr. 2018).
53. Xue, Y. *et al.* Single-Photon Emission from Point Defects in Aluminum Nitride Films. *The Journal of Physical Chemistry Letters* **11**, 2689–2694 (Mar. 2020).
54. Varley, J. B., Janotti, A. & de Walle, C. G. V. Defects in AlN as candidates for solid-state qubits. *Physical Review B* **93**. doi:[10.1103/physrevb.93.161201](https://doi.org/10.1103/physrevb.93.161201) (Apr. 2016).
55. Hardy, W. J. *et al.* Single and double hole quantum dots in strained Ge/SiGe quantum wells. *Nanotechnology* **30**, 215202 (Mar. 2019).
56. Toth, M. & Aharonovich, I. Single Photon Sources in Atomically Thin Materials. *Annual Review of Physical Chemistry* **70**, 123–142 (June 2019).
57. Atatüre, M., Englund, D., Vamivakas, N., Lee, S.-Y. & Wrachtrup, J. Material platforms for spin-based photonic quantum technologies. *Nature Reviews Materials* **3**, 38–51 (Apr. 2018).

58. Tran, T. T. *et al.* Robust Multicolor Single Photon Emission from Point Defects in Hexagonal Boron Nitride. *ACS Nano* **10**, 7331–7338 (July 2016).
59. Tran, T. T., Bray, K., Ford, M. J., Toth, M. & Aharonovich, I. *Quantum Emission from Hexagonal Boron Nitride Monolayers* in *Conference on Lasers and Electro-Optics* (OSA, 2016). doi:[10.1364/cleo_qels.2016.ftu4d.1](https://doi.org/10.1364/cleo_qels.2016.ftu4d.1).
60. Weston, L., Wickramaratne, D., Mackoït, M., Alkauskas, A. & de Walle, C. G. V. Native point defects and impurities in hexagonal boron nitride. *Physical Review B* **97**. doi:[10.1103/physrevb.97.214104](https://doi.org/10.1103/physrevb.97.214104) (June 2018).
61. Abdi, M., Chou, J.-P., Gali, A. & Plenio, M. B. Color Centers in Hexagonal Boron Nitride Monolayers: A Group Theory and Ab Initio Analysis. *ACS Photonics* **5**, 1967–1976 (Apr. 2018).
62. Persson, C. *Brief Introduction to the density functional theory* 2020.
63. Top500. SUPERCOMPUTER FUGAKU June 2020. <<https://www.top500.org/system/179807/>> (visited on 10/02/2020).
64. Hohenberg, P. & Kohn, W. Inhomogeneous Electron Gas. *Physical Review* **136**, B864–B871 (Nov. 1964).
65. Toulouse, J. *Introduction to density-functional theory* Sept. 2019. <http://www.lct.jussieu.fr/pagesperso/toulouse/enseignement/introduction_dft.pdf> (visited on 10/25/2020).
66. Kohn, W. & Sham, L. J. Self-Consistent Equations Including Exchange and Correlation Effects. *Physical Review* **140**, A1133–A1138 (Nov. 1965).
67. Allen, J. P. & Watson, G. W. Occupation matrix control of d- and f-electron localisations using DFT+U. *Phys. Chem. Chem. Phys.* **16**, 21016–21031 (2014).
68. David Sholl, J. A. S. *Density Functional Theory: A Practical Introduction* 238 pp. ISBN: 0470373172. <https://www.ebook.de/de/product/7207845/david_sholl_janice_a_steckel_density_functional_theory_a_practical_introduction.html> (WILEY, 2009).
69. Perdew, J. P. & Wang, Y. Accurate and simple analytic representation of the electron-gas correlation energy. *Physical Review B* **45**, 13244–13249 (June 1992).
70. Perdew, J. P., Burke, K. & Ernzerhof, M. Generalized Gradient Approximation Made Simple. *Physical Review Letters* **77**, 3865–3868 (Oct. 1996).
71. Freysoldt, C. *et al.* First-principles calculations for point defects in solids. *Reviews of Modern Physics* **86**, 253–305 (Mar. 2014).

72. Tran, F. & Blaha, P. Accurate Band Gaps of Semiconductors and Insulators with a Semilocal Exchange–Correlation Potential. *Physical Review Letters* **102**. doi:[10.1103/physrevlett.102.226401](https://doi.org/10.1103/physrevlett.102.226401) (June 2009).
73. Becke, A. D. & Johnson, E. R. A simple effective potential for exchange. *The Journal of Chemical Physics* **124**, 221101 (June 2006).
74. Koller, D., Tran, F. & Blaha, P. Merits and limits of the modified Becke–Johnson exchange potential. *Physical Review B* **83**. doi:[10.1103/physrevb.83.195134](https://doi.org/10.1103/physrevb.83.195134) (May 2011).
75. Becke, A. D. A new mixing of Hartree–Fock and local density-functional theories. *The Journal of Chemical Physics* **98**, 1372–1377 (Jan. 1993).
76. Perdew, J. P., Ernzerhof, M. & Burke, K. Rationale for mixing exact exchange with density functional approximations. *The Journal of Chemical Physics* **105**, 9982–9985 (Dec. 1996).
77. Aryasetiawan, F. & Gunnarsson, O. The GW method. *Reports on Progress in Physics* **61**, 237–312 (Mar. 1998).
78. Heyd, J., Scuseria, G. E. & Ernzerhof, M. Hybrid functionals based on a screened Coulomb potential. *The Journal of Chemical Physics* **118**, 8207–8215 (May 2003).
79. Krukau, A. V., Vydrov, O. A., Izmaylov, A. F. & Scuseria, G. E. Influence of the exchange screening parameter on the performance of screened hybrid functionals. *The Journal of Chemical Physics* **125**, 224106 (Dec. 2006).
80. Klimeš, J., Bowler, D. R. & Michaelides, A. Chemical accuracy for the van der Waals density functional. *Journal of Physics: Condensed Matter* **22**, 022201 (Dec. 2009).
81. Dion, M., Rydberg, H., Schröder, E., Langreth, D. C. & Lundqvist, B. I. Van der Waals Density Functional for General Geometries. *Physical Review Letters* **92**. doi:[10.1103/physrevlett.92.246401](https://doi.org/10.1103/physrevlett.92.246401) (June 2004).
82. Freitas, L. C. G. Prêmio Nobel de Química em 1998: Walter Kohn e John A. Pople. *Química Nova* **22**, 293–298 (Apr. 1999).
83. Murphy, K. *Machine learning : a probabilistic perspective* ISBN: 9780262018029 (MIT Press, Cambridge, Mass, 2012).
84. Wolpert, D. & Macready, W. No Free Lunch Theorems for Search (Mar. 1996).
85. Guido, S. *Introduction to Machine Learning with Python* 400 pp. ISBN: 1449369413. <https://www.ebook.de/de/product/23308778/sarah_guido_introduction_to_machine_learning_with_python.html> (O’Reilly UK Ltd., 2016).

86. Caruana, R. & Niculescu-Mizil, A. *An empirical comparison of supervised learning algorithms* in *Proceedings of the 23rd international conference on Machine learning - ICML '06* (ACM Press, 2006). doi:[10.1145/1143844.1143865](https://doi.org/10.1145/1143844.1143865).
87. Friedman, J., Hastie, T. & Tibshirani, R. Additive logistic regression: a statistical view of boosting (With discussion and a rejoinder by the authors). *The Annals of Statistics* **28**, 337–407 (Apr. 2000).
88. Friedman, J. H. Greedy function approximation: A gradient boosting machine. *Ann. Statist.* **29**, 1189–1232 (Oct. 2001).
89. James, G., Witten, D., Hastie, T. & Tibshirani, R. *An Introduction to Statistical Learning* ISBN: 1461471370. <https://www.ebook.de/de/product/20292548/gareth_james_daniela_witten_trevor_hastie_robert_tibshirani_an_introduction_to_statistical_learning.html> (Springer-Verlag GmbH, 2017).
90. Marsland, S. *Machine Learning* 457 pp. ISBN: 9781466583337 (Taylor & Francis Ltd., 2014).
91. Kresse, G. & Furthmüller, J. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. *Physical Review B* **54**, 11169–11186 (Oct. 1996).
92. Battle, R. & Benson, E. Bridging the semantic Web and Web 2.0 with Representational State Transfer (REST). *Journal of Web Semantics* **6**, 61–69 (Feb. 2008).
93. Kirklin, S. *et al.* The Open Quantum Materials Database (OQMD): assessing the accuracy of DFT formation energies. *npj Computational Materials* **1**. doi:[10.1038/npjcompumats.2015.10](https://doi.org/10.1038/npjcompumats.2015.10) (Dec. 2015).
94. Saal, J. E., Kirklin, S., Aykol, M., Meredig, B. & Wolverton, C. Materials Design and Discovery with High-Throughput Density Functional Theory: The Open Quantum Materials Database (OQMD). *JOM* **65**, 1501–1509 (Sept. 2013).
95. Jain, A. *et al.* Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Materials* **1**, 011002 (July 2013).
96. Ong, S. P. *et al.* The Materials Application Programming Interface (API): A simple, flexible and efficient API for materials data based on Representational State Transfer (REST) principles. *Computational Materials Science* **97**, 209–215 (Feb. 2015).
97. Levin, I. *NIST Inorganic Crystal Structure Database (ICSD)* en. 2020. doi:[10.18434/M32147](https://doi.org/10.18434/M32147).

98. Landis, D. D. *et al.* The Computational Materials Repository. *Computing in Science & Engineering* **14**, 51–57 (Nov. 2012).
99. Larsen, A. H. *et al.* The atomic simulation environment—a Python library for working with atoms. *Journal of Physics: Condensed Matter* **29**, 273002 (June 2017).
100. Ong, S. P. *et al.* Python Materials Genomics (pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science* **68**, 314–319 (Feb. 2013).
101. Ward, L. *et al.* Matminer: An open source toolkit for materials data mining. *Computational Materials Science* **152**, 60–69 (Sept. 2018).
102. Draxl, C. & Scheffler, M. The NOMAD laboratory: from data sharing to artificial intelligence. *Journal of Physics: Materials* **2**, 036001 (May 2019).
103. Wang, L., Maxisch, T. & Ceder, G. Oxidation energies of transition metal oxides within the GGA+U framework. *Physical Review B* **73**. doi:[10.1103/physrevb.73.195107](https://doi.org/10.1103/physrevb.73.195107) (May 2006).
104. Curtarolo, S. *et al.* AFLOWLIB.ORG: A distributed materials properties repository from high-throughput ab initio calculations. *Computational Materials Science* **58**, 227–235 (June 2012).
105. Curtarolo, S. *et al.* AFLOW: An automatic framework for high-throughput materials discovery. *Computational Materials Science* **58**, 218–226 (June 2012).
106. Calderon, C. E. *et al.* The AFLOW standard for high-throughput materials science calculations. *Computational Materials Science* **108**, 233–238 (Oct. 2015).
107. Setyawan, W. & Curtarolo, S. High-throughput electronic band structure calculations: Challenges and tools. *Computational Materials Science* **49**, 299–312 (Aug. 2010).
108. Setyawan, W., Gaume, R. M., Lam, S., Feigelson, R. S. & Curtarolo, S. High-Throughput Combinatorial Database of Electronic Band Structures for Inorganic Scintillator Materials. *ACS Combinatorial Science* **13**, 382–390 (June 2011).
109. Isayev, O. *et al.* Universal fragment descriptors for predicting properties of inorganic crystals. *Nature Communications* **8**. doi:[10.1038/ncomms15679](https://doi.org/10.1038/ncomms15679) (June 2017).
110. Ferrenti, A. M., de Leon, N. P., Thompson, J. D. & Cava, R. J. Identifying candidate hosts for quantum defects via data mining. *npj Computational Materials* **6**. doi:[10.1038/s41524-020-00391-7](https://doi.org/10.1038/s41524-020-00391-7) (Aug. 2020).

111. Stevanović, V., Lany, S., Zhang, X. & Zunger, A. Correcting density functional theory for accurate predictions of compound enthalpies of formation: Fitted elemental-phase reference energies. *Physical Review B* **85**. doi:[10.1103/physrevb.85.115104](https://doi.org/10.1103/physrevb.85.115104) (Mar. 2012).
112. Choudhary, K. *et al.* JARVIS: An Integrated Infrastructure for Data-driven Materials Design. arXiv: [2007.01831v1](https://arxiv.org/abs/2007.01831) [[cond-mat.mtrl-sci](#)] (July 3, 2020).
113. Thonhauser, T. *et al.* Van der Waals density functional: Self-consistent potential and the nature of the van der Waals bond. *Physical Review B* **76**. doi:[10.1103/physrevb.76.125112](https://doi.org/10.1103/physrevb.76.125112) (Sept. 2007).
114. Klimeš, J., Bowler, D. R. & Michaelides, A. Van der Waals density functionals applied to solids. *Physical Review B* **83**. doi:[10.1103/physrevb.83.195131](https://doi.org/10.1103/physrevb.83.195131) (May 2011).
115. Choudhary, K., Cheon, G., Reed, E. & Tavazza, F. Elastic properties of bulk and low-dimensional materials using van der Waals density functional. *Physical Review B* **98**. doi:[10.1103/physrevb.98.014107](https://doi.org/10.1103/physrevb.98.014107) (July 2018).
116. Choudhary, K. *et al.* Computational screening of high-performance optoelectronic materials using OptB88vdW and TB-mBJ formalisms. *Scientific Data* **5**. doi:[10.1038/sdata.2018.82](https://doi.org/10.1038/sdata.2018.82) (May 2018).
117. Ohebbi. *ohebbi/predicting-solid-state-qubit-candidates: v0.1-beta* 2021. doi:[10.5281/ZENODO.4633959](https://doi.org/10.5281/ZENODO.4633959).
118. Rosenbrock, C. W. A Practical Python API for Querying AFLOWLIB. arXiv: [1710.00813v1](https://arxiv.org/abs/1710.00813) [[cs.DB](#)] (Sept. 28, 2017).
119. Breuck, P.-P. D., Evans, M. L. & Rignanese, G.-M. Robust model benchmarking and bias-imbalance in data-driven materials science: a case study on MODNet. arXiv: [2102.02263v1](https://arxiv.org/abs/2102.02263) [[cond-mat.mtrl-sci](#)] (Feb. 3, 2021).
120. Markham, M. *et al.* CVD diamond for spintronics. *Diamond and Related Materials* **20**, 134–139 (Feb. 2011).
121. Balasubramanian, G. *et al.* Ultralong spin coherence time in isotopically engineered diamond. *Nature Materials* **8**, 383–387 (Apr. 2009).
122. Tyryshkin, A. M. *et al.* Electron spin coherence exceeding seconds in high-purity silicon. *Nature Materials* **11**, 143–147 (Dec. 2011).
123. Inselberg, A. The plane with parallel coordinates. *The Visual Computer* **1**, 69–91 (Aug. 1985).

124. Inselberg, A. & Dimsdale, B. *Parallel coordinates: a tool for visualizing multi-dimensional geometry* in *Proceedings of the First IEEE Conference on Visualization: Visualization 90* (IEEE Comput. Soc. Press, 1990). doi:[10.1109/visual.1990.146402](https://doi.org/10.1109/visual.1990.146402).
125. Ericson, D., Johansson, J. & Cooper, M. *Visual Data Analysis using Tracked Statistical Measures within Parallel Coordinate Representations in Coordinated and Multiple Views in Exploratory Visualization (CMV'05)* (IEEE). doi:[10.1109/cmv.2005.21](https://doi.org/10.1109/cmv.2005.21).
126. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* (2011). arXiv: [1201.0490v4](https://arxiv.org/abs/1201.0490v4) [[cs.LG](#)] (Jan. 2, 2012).
127. Lemaitre, G., Nogueira, F. & Aridas, C. K. Imbalanced-learn: A Python Toolbox to Tackle the Curse of Imbalanced Datasets in Machine Learning. arXiv: [1609.06570v1](https://arxiv.org/abs/1609.06570v1) [[cs.LG](#)] (Sept. 21, 2016).
128. Kotochigova, S., Levine, Z. H., Shirley, E. L., Stiles, M. D. & Clark, C. W. Local-density-functional calculations of the energy of atoms. *Physical Review A* **55**, 191–199 (Jan. 1997).
129. Laws, K. J., Miracle, D. B. & Ferry, M. A predictive structural model for bulk metallic glasses. *Nature Communications* **6**. doi:[10.1038/ncomms9123](https://doi.org/10.1038/ncomms9123) (Sept. 2015).
130. Butler, M. A. & Ginley, D. S. Prediction of Flatband Potentials at Semiconductor-Electrolyte Interfaces from Atomic Electronegativities. *Journal of The Electrochemical Society* **125**, 228–232 (Feb. 1978).
131. Ward, L., Agrawal, A., Choudhary, A. & Wolverton, C. A general-purpose machine learning framework for predicting properties of inorganic materials. *npj Computational Materials* **2**. doi:[10.1038/npjcompumats.2016.28](https://doi.org/10.1038/npjcompumats.2016.28) (Aug. 2016).
132. Deml, A. M., O'Hayre, R., Wolverton, C. & Stevanović, V. Predicting density functional theory total energies and enthalpies of formation of metal-nonmetal compounds by linear regression. *Physical Review B* **93**. doi:[10.1103/physrevb.93.085142](https://doi.org/10.1103/physrevb.93.085142) (Feb. 2016).
133. Weeber, A. W. Application of the Miedema model to formation enthalpies and crystallisation temperatures of amorphous alloys. *Journal of Physics F: Metal Physics* **17**, 809–813 (Apr. 1987).
134. Yang, X. & Zhang, Y. Prediction of high-entropy stabilized solid-solution in multi-component alloys. *Materials Chemistry and Physics* **132**, 233–238 (Feb. 2012).

135. Rupp, M., Tkatchenko, A., Müller, K.-R. & von Lilienfeld, O. A. Fast and Accurate Modeling of Molecular Atomization Energies with Machine Learning. *Physical Review Letters* **108**. doi:[10.1103/physrevlett.108.058301](https://doi.org/10.1103/physrevlett.108.058301) (Jan. 2012).
136. Schütt, K. T. *et al.* How to represent crystal structures for machine learning: Towards fast prediction of electronic properties. *Physical Review B* **89**. doi:[10.1103/physrevb.89.205118](https://doi.org/10.1103/physrevb.89.205118) (May 2014).
137. Faber, F., Lindmaa, A., von Lilienfeld, O. A. & Armiento, R. Crystal structure representations for machine learning models of formation energies. *International Journal of Quantum Chemistry* **115**, 1094–1101 (Apr. 2015).
138. Ewald, P. P. Die Berechnung optischer und elektrostatischer Gitterpotentiale. *Annalen der Physik* **369**, 253–287 (1921).
139. Hansen, K. *et al.* Machine Learning Predictions of Molecular Properties: Accurate Many-Body Potentials and Nonlocality in Chemical Space. *The Journal of Physical Chemistry Letters* **6**, 2326–2331 (June 2015).
140. Ward, L. *et al.* Including crystal structure attributes in machine learning models of formation energies via Voronoi tessellations. *Physical Review B* **96**. doi:[10.1103/physrevb.96.024104](https://doi.org/10.1103/physrevb.96.024104) (July 2017).
141. Botu, V. & Ramprasad, R. Adaptive machine learning framework to accelerate ab initio molecular dynamics. *International Journal of Quantum Chemistry* **115**, 1074–1083 (Dec. 2014).
142. De Jong, M. *et al.* A Statistical Learning Framework for Materials Science: Application to Elastic Moduli of k-nary Inorganic Polycrystalline Compounds. *Scientific Reports* **6**. doi:[10.1038/srep34256](https://doi.org/10.1038/srep34256) (Oct. 2016).
143. Seko, A., Hayashi, H., Nakayama, K., Takahashi, A. & Tanaka, I. Representation of compounds for machine-learning prediction of physical properties. *Physical Review B* **95**. doi:[10.1103/physrevb.95.144110](https://doi.org/10.1103/physrevb.95.144110) (Apr. 2017).
144. Steinhardt, P. J., Nelson, D. R. & Ronchetti, M. Bond-orientational order in liquids and glasses. *Physical Review B* **28**, 784–805 (July 1983).
145. Waroquiers, D. *et al.* Statistical Analysis of Coordination Environments in Oxides. *Chemistry of Materials* **29**, 8346–8360 (Sept. 2017).
146. Zimmermann, N. E. R., Horton, M. K., Jain, A. & Haranczyk, M. Assessing Local Structure Motifs Using Order Parameters for Motif Recognition, Interstitial Identification, and Diffusion Path Characterization. *Frontiers in Materials* **4**. doi:[10.3389/fmats.2017.00034](https://doi.org/10.3389/fmats.2017.00034) (Nov. 2017).

147. Behler, J. Atom-centered symmetry functions for constructing high-dimensional neural network potentials. *The Journal of Chemical Physics* **134**, 074106 (Feb. 2011).
148. Khorshidi, A. & Peterson, A. A. Amp: A modular approach to machine learning in atomistic simulations. *Computer Physics Communications* **207**, 310–324 (Oct. 2016).
149. Peng, H. L., Li, M. Z. & Wang, W. H. Structural Signature of Plastic Deformation in Metallic Glasses. *Physical Review Letters* **106**. doi:[10.1103/physrevlett.106.135503](https://doi.org/10.1103/physrevlett.106.135503) (Mar. 2011).
150. Wang, Q. & Jain, A. A transferable machine-learning framework linking interstice distribution and plastic heterogeneity in metallic glasses. *Nature Communications* **10**. doi:[10.1038/s41467-019-13511-9](https://doi.org/10.1038/s41467-019-13511-9) (Dec. 2019).
151. Dylla, M. T., Dunn, A., Anand, S., Jain, A. & Snyder, G. J. Machine Learning Chemical Guidelines for Engineering Electronic Structures in Half-Heusler Thermoelectric Materials. *Research* **2020**, 1–8 (Apr. 2020).
152. None Available. *Materials Data on PH6C2S2N(ClO2)2 by Materials Project* en. 2020. doi:[10.17188/1268877](https://doi.org/10.17188/1268877).
153. None Available. *Materials Data on Nb7S2I19 by Materials Project* en. 2014. doi:[10.17188/1277059](https://doi.org/10.17188/1277059).