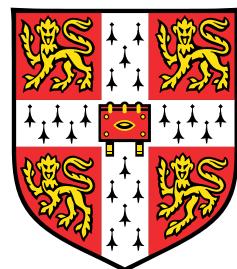


Repurposing Drugs for the Rapid Response to Epidemics and Pandemics

Using Batch Active Learning



Ross Brown

Department of Chemical Engineering and Biotechnology
University of Cambridge

This dissertation is submitted for the degree of
Master of Engineering

Robinson College

May 2022

Declaration

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other university. This dissertation is my own work and contains nothing which is the outcome of work done in collaboration with others, except as specified in the text and Acknowledgements. This dissertation contains fewer than 10,000 words including appendices, bibliography, footnotes, tables and equations and has fewer than 40 pages.

Ross Brown
May 2022

Acknowledgements

And I would like to acknowledge ...

Abstract

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Table of contents

List of figures	xi
Nomenclature	xiii
1 Introduction	1
2 Previous Work	3
2.1 Active Learning	4
2.1.1 Current Data	5
2.1.2 Estimated Future	11
2.2 Batch Active Learning	11
2.3 Drug Data for Machine Learning	14
2.3.1 Physical Properties	14
2.3.2 Fingerprints	15
3 Methodology	17
3.1 Outline	17
3.2 Proof	18
3.2.1 Custom Base Functions	18
3.2.2 Active Learning Algorithms	20
3.2.3 Training Framework	20
4 Results	23
4.1 Overview	23
4.2 Second Part	24
4.3 Special Case: COVID-19	24
5 Discussion	25

6 Conclusion	27
References	29

List of figures

2.1	Example Dataset for Representation of Ideas	4
2.2	Uncertainty Sampling Demonstration	6
2.3	Broad-Base Sampling Illustration	9
2.5	Batch Uncertainty Sampling	12
2.6	Batch Broad-Base Sampling	13
3.1	Graphical representation of the split function. The red dashed boundary represents the input (additional colour coding has been performed to assist the reader in understanding the transposition of the base components).	19

Nomenclature

Chapter 2

N	Number of features/dimensions of x
x	Data points where $x = \{x_0, x_1, \dots, x_{N-1}\}$
y	Labels for the dataset where $y = \{y_0, y_1, \dots, y_{N-1}\}$

Chapter 3

X_{test}	Datasets used to provide a score for the algorithms
X_{train}	Datasets used for training the algorithms
x_{known}	Data points where the true label is available to the algorithms used
x_{unknown}	Data points where the true label is not available to the algorithms used
y_{known}	True labels available to the algorithms used
y_{unknown}	True labels unavailable to the algorithms used

Chapter 1

Introduction

In 2019, human civilisation was on the precipice of a natural disaster: SARS-CoV-2 (COVID-19). First reported to the WHO on December 31st, it became officially recognised as a pandemic on March 11th 2020. As of the writing of this passage, 515 million cases and 15 million excess deaths have been recorded. This, however, is not the first time a pandemic has occurred, with the Black Death infamously killing a third of Europe's population and the Spanish Flu causing mass death throughout the world. Likewise, it is unlikely to be the last.

When such a disaster does strike, it is important to react quickly. Vaccinations are allowed accelerated timelines in development cutting development from years to month, and trials into potential treatments are encouraged with haste. Within the first stages of the pandemic, drugs such as hydroxychloroquine and bleach were amongst several that were promoted by the President of the United States of America demonstrating the desperation in finding therapeutic drugs against the virus.

In order to facilitate a more robust approach to finding treatments, the FDA instigated the Coronavirus Treatment Acceleration Program (CTAP) [Cen22]. Here, over 690 drugs are in the development stage with over 450 clinical trials underway to investigate the effectiveness, with 15 drugs currently authorised for emergency use and only one drug, remdesivir, with approval for use against COVID-19 [Cen22]. Indeed, remdesivir is an important case. This drug was developed initially for hepatitis-c before being used for several other conditions until finally being used for COVID-19 [Par+20]. This demonstrates how a discovered drug can be repurposed for new diseases providing a cheap means of drug "redevelopment".

Investigations into pre-existing drugs, however, were slow and largely carried out through labour intensive mechanisms without a rational methodical testing regime. This added time to finding treatments to COVID-19. Time many did not have. A hopeful fulfilment of this problem is the "Robot Scientist"; a fully automated combination of software and hardware

aimed at solving this problem. For the software side, a form of reinforcement machine learning is proposed: active learning. This is a methodology suited to fields with large amounts of unlabelled data which is difficult to label. In this case, the labelling requires chemical and biological experimentation costing both time and money. By using Active Learning, as few drugs as possible will be labelled within this stage to accurately predict the best drugs for the given problem. From here, clinical trials may begin. Additionally, due to the large importance of time, many drugs may be tested in parallel. This presents an additional problem: how does one set up a testing scheme for batches.

Thus, the purpose of this thesis. To present an algorithm which may be used to discover effective drugs within a short period of time. Additionally, a framework will be developed that allows for different algorithms to be rigorously compared to each other for increased robustness.

Chapter 2

Previous Work

Scores displayed in examples have been based on the entire data set. Although this usually leads to data leakage within machine learning, this is not a concern here as the true comparison comes from testing *intelligent* vs *dumb* learning methods. In both of these cases, the model is kept identical, but the selection process is not. The baseline simply takes the first n entries from the data set, with the *intelligent* method described where required. A simple function has been used to present data as a means of demonstrating these models, with x having two dimensions. The function for y is shown in 2.1 and displayed graphically in Figure 2.1.

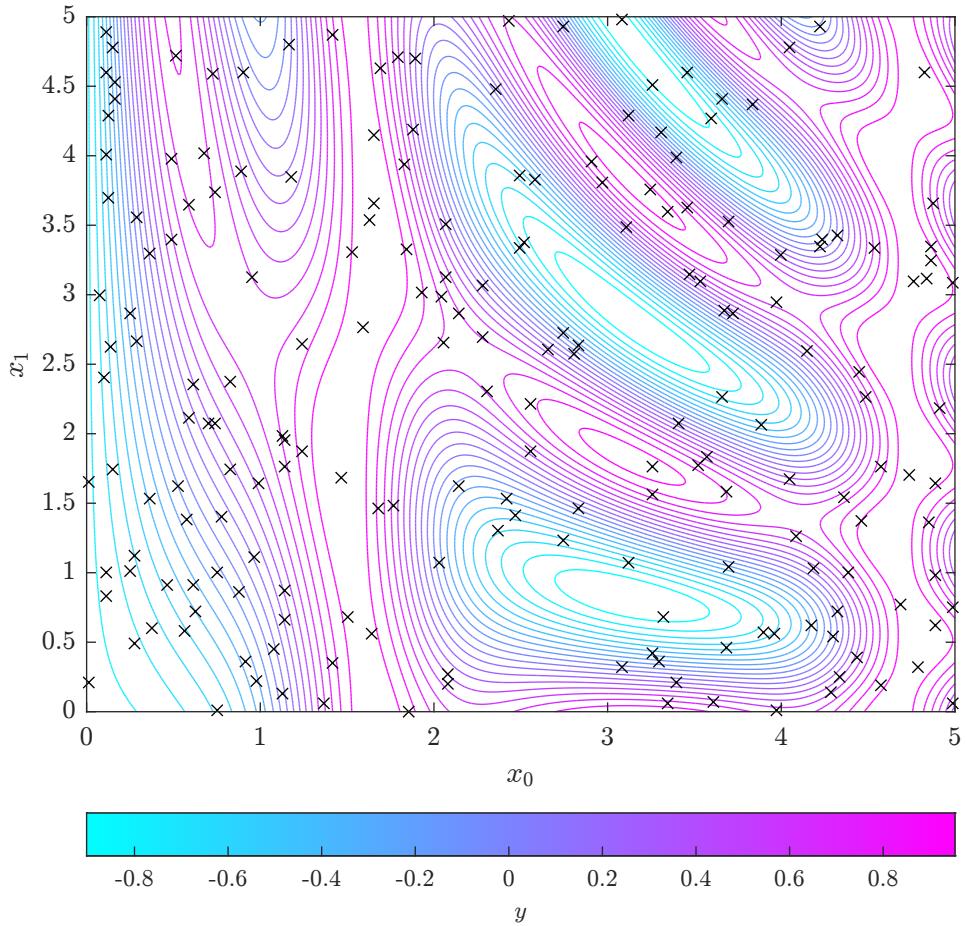


Fig. 2.1 Contour plot of the function used to demonstrate the algorithms presented in previous work. The crosses have been used to show the location of the 200 test data points used within this example.

$$y = \sin(x_1)^{10} + \cos(10 + x_1 x_2) \cos(x_1) \quad (2.1)$$

2.1 Active Learning

There are several schools of thought regarding active learning. These can be separated into two distinct categories: current data and future predictions. The former of these is computationally cheaper, as will be apparent on description.

2.1.1 Current Data

Uncertainty Sampling and Regions of Disagreements

The simplest is applicable to cases in which a certainty is provided with each prediction. Settles [Set09] suggests selecting the data point with the largest uncertainty according to the current model. Using the dataset ”, this is demonstrated in Figure 2.2 with the algorithm for deciding the next sample point given in Algorithm 1.

Algorithm 1: Uncertainty Sampling Selection

Data: $X_{\text{known}}, Y_{\text{known}}, X_{\text{unknown}}$
Result: Next X to label
model = BayesianRidge();
model.fit($X_{\text{known}}, Y_{\text{known}}$);
standard_deviation = model.standard_deviation(X_{unknown});
return $\max(\text{standard_deviation})$

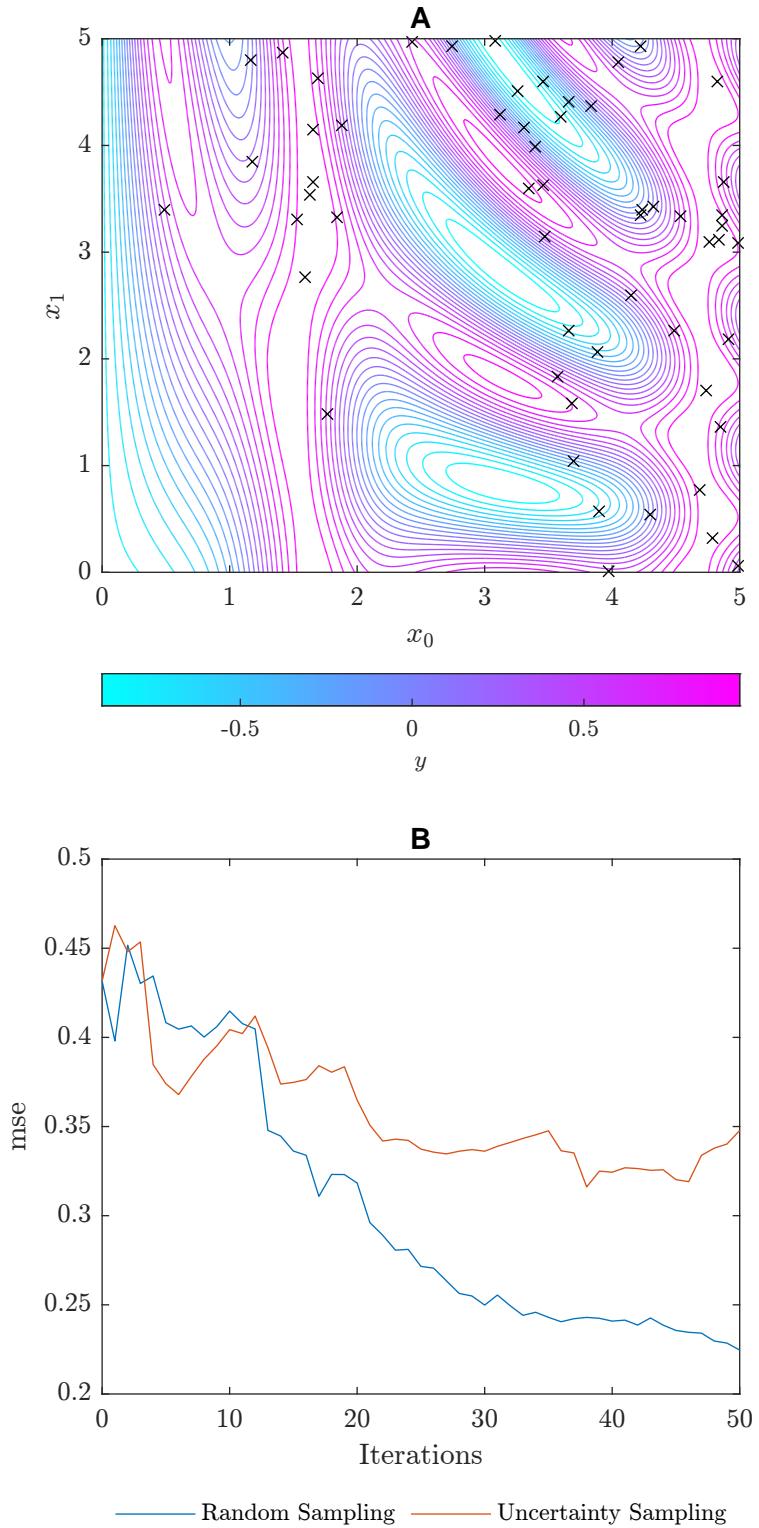


Fig. 2.2 The outcome of the investigating the areas of the highest uncertainty. An initial set of 5 random points was provided, and 50 further iterations were then carried out of sample size 1. A) Demonstrates the final set of points tested by the algorithm and B) shows the change in the mean squared error for the algorithm after each iteration.

Interestingly, Figure 2.2B shows how the mean squared error for the random sampling method performed to worse within the iterations tested. This is likely due to a bias in the use of linear models in fitting leading to large uncertainties surrounding areas with high curvature. Evidence to this is provided in 2.2A with a large proportion of the sampled points at areas of high curvature.

As addressed by Settles [Set09], this can be extended to any probabilistic model through 2.2. Settles [Set09] also notes the use of information theory for probabilistic models(2.3), where y_i refers to all possible categorisations for x . This derives from the principle that the greatest entropy requires the most information to encode, and thus the least certain. However, Settles [Set09] fails to address non-probabilistic models in this instance, instead converting such models into probabilistic ones.

$$x_{\text{next}} = \underset{X}{\operatorname{argmax}} [s_{g(X)}] \quad (2.2)$$

In order to adapt non-probabilistic models into probabilistic ones, composite models may be used. These are an amalgamation of other models where the standard deviation of the individual models can be taken as the degree of certainty for a given point. Many authors have called this as minimising the region of disagreement as it attempts to produce a coherent hypothesis space. By minimising the region of disagreement between various models, a finer fit may be achieved. Indeed, this was the method used in Figure 2.2.

One way of achieving this, especially in a regression model where boundaries are not quite so distinct, is to declare n models $M = \{m_1, \dots, m_n\}$. Combining these allow for a model \hat{m} to be defined with prediction \hat{y} , being the mean prediction of M , $\frac{1}{n} \sum y_i$ and a sample standard deviation \hat{s} defined as the sample standard deviation of y_i . This standard deviation can be used as a measure of the disagreement between the models. Thus, using a method as in Section ??.

$$x_{\text{next}} = \underset{x}{\operatorname{argmax}} \left[- \sum_i P(y_i|x) \ln P(y_i|x) \right] \quad (2.3)$$

Broad Knowledge Base

A second form stems from information theory. Here, the aim is to produce an evenly dispersed x allowing a well-informed knowledge base. This prevents poor model choice from influencing the algorithm as was seen in 2.2. There are two paths to proceed: density and nearest neighbours.

The former of these requires a definition of density in a sparsely populated space. As an analogy, although the density of a gas appears well-defined, it becomes non-smooth once the

volume defined over is comparable to the distance between particles. Thus, a new definition is required.

Alternatively, nearest neighbour requires little explanation. x_{next} is the unlabelled data point furthest from any labelled data point.

$$x_{\text{next}} = \operatorname{argmax}_x \left(\sum \frac{1}{\text{sim}(x, x_i)} \right) \quad (2.4)$$

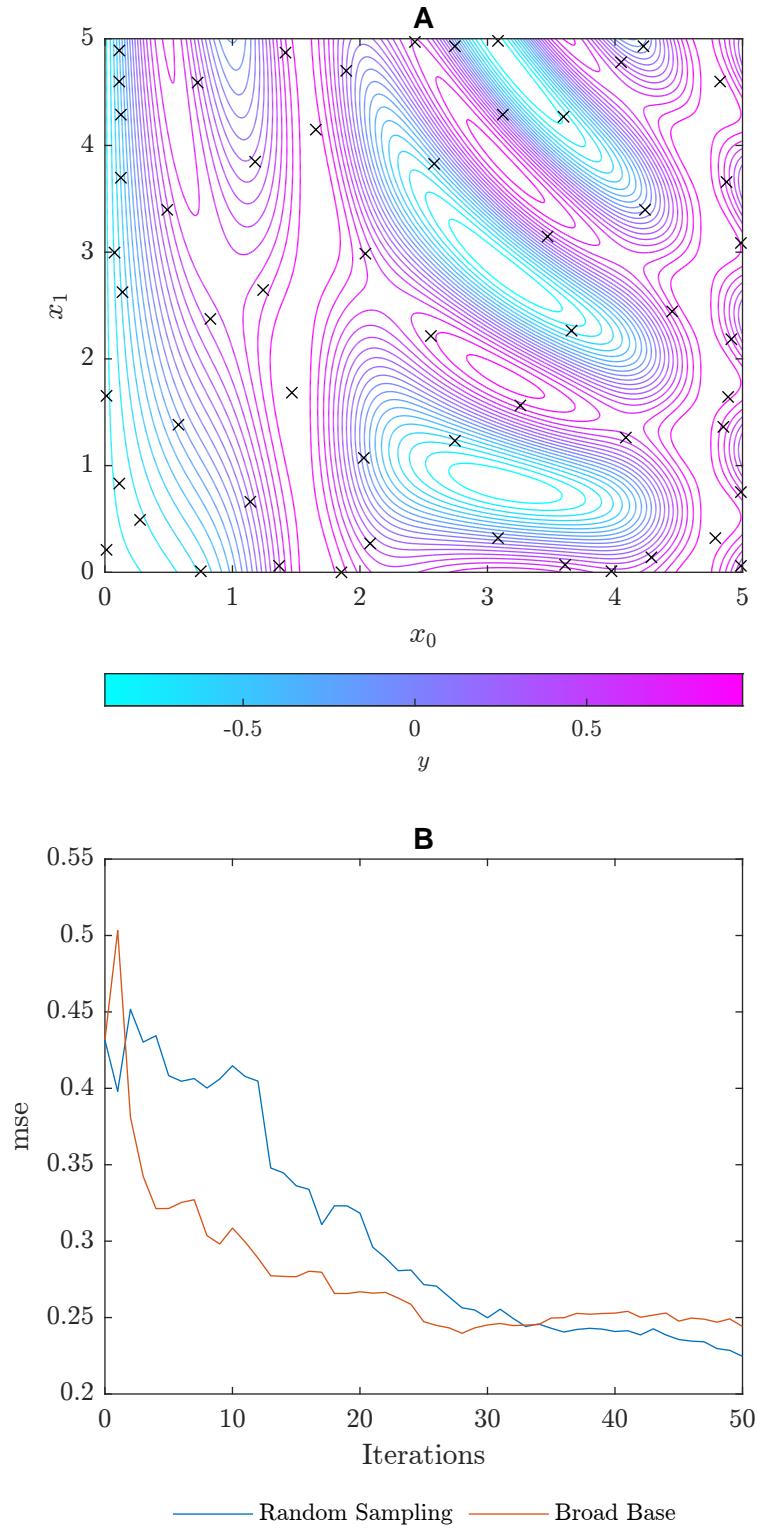


Fig. 2.3 The outcome of the investigating the areas of using a broad base. An initial set of 5 random points was provided, and 50 further iterations were then carried out of sample size 1. A) Demonstrates the final set of points tested by the algorithm and B) shows the change in the mean squared error for the algorithm after each iteration.

Density Hotspots

Conversely, a density weighted model has been suggested, as it escapes the introduction of error from outliers (i.e. data points far away from alternative data points). Settles and Craven [SC08] suggest (2.5) which can be broken down into two parts: a function for selection, ϕ_A , and a function for similarity, sim. The former arises from another method described in this section. The latter requires a function to describe the similarity between data points.

$$x_{\text{next}} = \underset{x}{\operatorname{argmax}} \left[\phi_A(x) \times \left(\frac{1}{U} \sum \text{sim}(x, x_i) \right)^\beta \right] \quad (2.5)$$

Settles and Craven [SC08] admits that sim is open for interpretation. It is also recognised that this lays the foundation of a clusterisation algorithm. There exist many forms of these algorithms, with the results of several of these algoritms on toy data sets presented in Figure 2.4 [Sci].

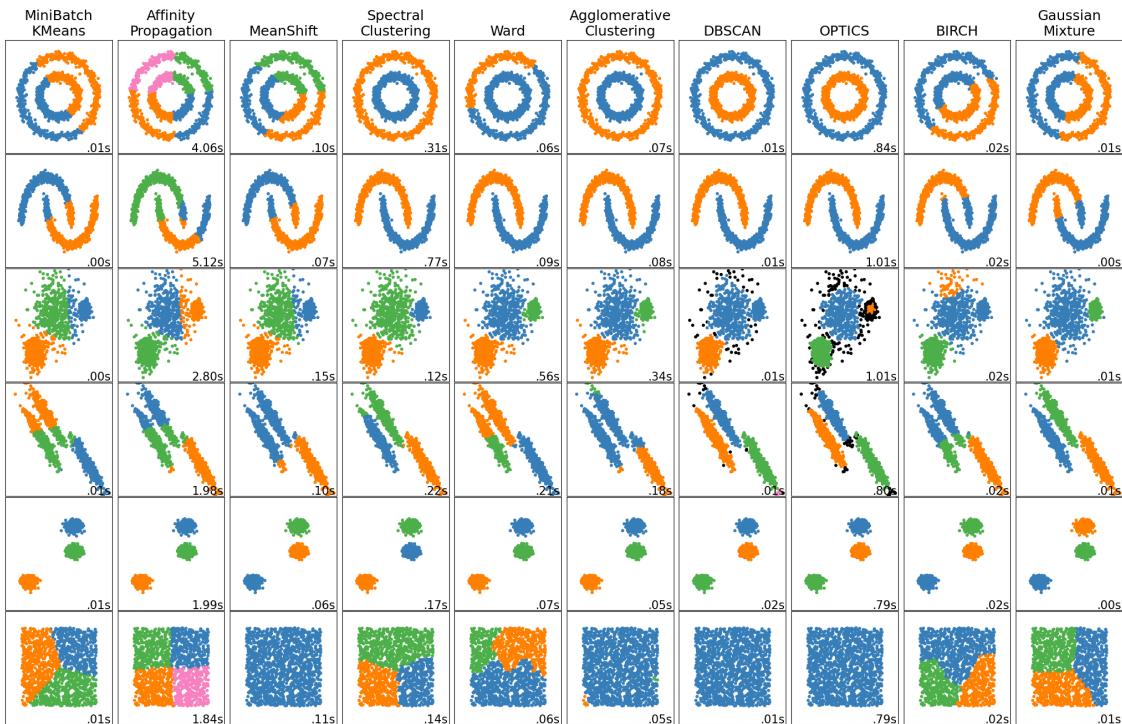


Fig. 2.4 Clusterisation algorithms used on sample two-dimensional data sets to demonstrate resultant clusters.

As demonstrated by 2.4 demonstrates, there are multiple different interpretations of the solution to the problem of clustering. The makers of the Sci-kit learn package also discuss the scalability of each algorithm [Sci]. In order to prepare a high number of features (beyond the

two used within this section for demonstration) and large number of data points, it is required that the algorithm scales accordingly. Further, for an adaptive process, it is more suitable for an algorithm to be adaptive to differing distribution. This limits the suitable algorithms to K-Means, Ward and Birch - columns one, five, and nine of Figure 2.4 respectively.

2.1.2 Estimated Future

These methods attempt to minimise a future attribute of the model. This works by predicting changes given with the inclusion of more data.

Expected Model Change

As the name implies, this method chooses points which are likely to have the largest impact on the final model. By instigating each potential point, the impact on the eventual model can be found. However, this requires a method for quantifying the model change.

Settles and Craven [SC08] and Settles [Set09] investigate models which can be trained "online": i.e. models which can use the previous iteration to reduce the time taken for convergence. They present a method called "Expected Gradient Length" (EGL) which has a couple of prerequisites: **1)** A probabilistic model is used **2)** Linear gradient based optimisation is used **3)** The model can be improved from previous iterations. Given these prerequisites, the problem becomes less computationally inexpensive given a small dataset or extensive parallelisation, and scales as $\mathcal{O}(n)$. However, it does have the distinct drawback of requiring close control of the data models used. Here, the length of the training gradient (the gradient used in re-fitting the parameters with gradient based optimisation) can be used as a measure of model change. In the case of a small model change, as is expected, the length of the training gradient can be written as $\|\nabla l(\langle x, y_i \rangle; \theta)\|$. Combining this with the probability distribution of y , the next sample to undergo labelling is given by 2.6.

$$x_{EGL}^* = \operatorname{argmax}_x P(y_i|x; \theta) \|\nabla l(\langle x, y_i \rangle; \theta)\| \quad (2.6)$$

2.2 Batch Active Learning

Several naive methods are available here. Firstly, getting the top N data points from a model described in Section. However, this method does not take into account the equivalence of the data points. This is extremely clear using the highest uncertainty method. Each method in Section[] has been modified to demonstrate this weakness.

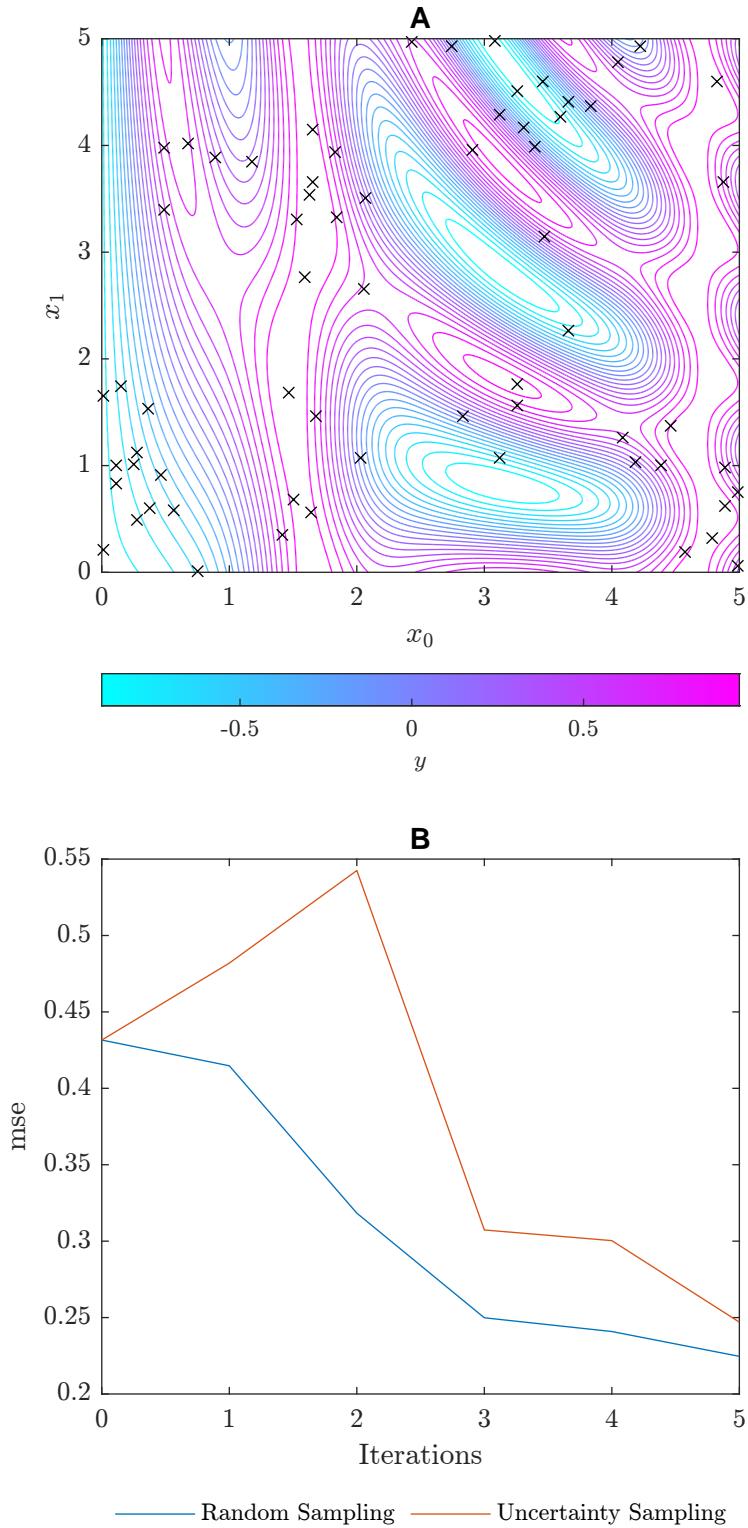


Fig. 2.5 The outcome of the investigating the areas of using uncertainty sampling. An initial set of 5 random points was provided, and 5 further iterations were then carried out of sample size 10. A) Demonstrates the final set of points tested by the algorithm and B) shows the change in the mean squared error for the algorithm after each iteration.

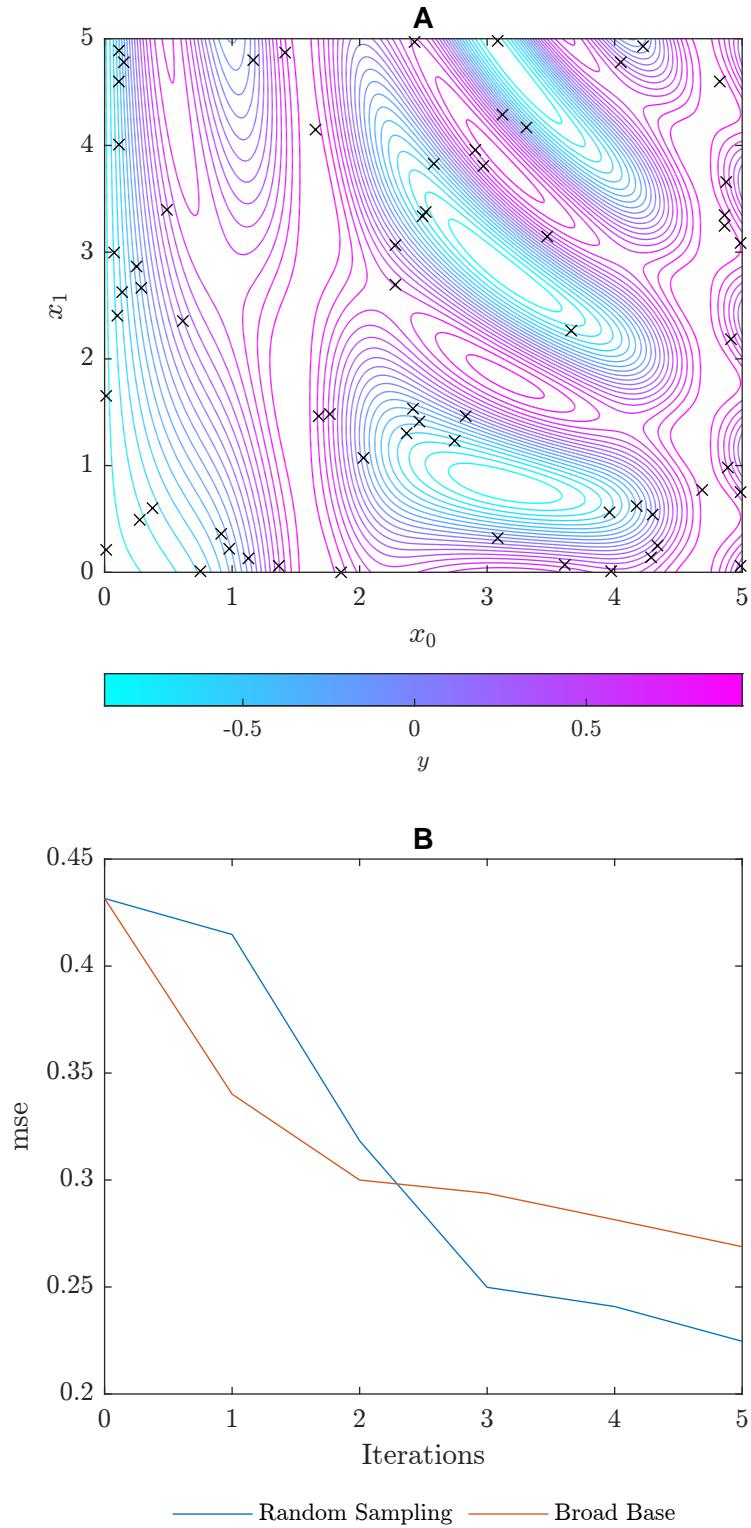


Fig. 2.6 The outcome of the investigating the areas of using broad-base sampling. An initial set of 5 random points was provided, and 5 further iterations were then carried out of sample size 10. A) Demonstrates the final set of points tested by the algorithm and B) shows the change in the mean squared error for the algorithm after each iteration.

It stands to reason that the area which has the highest uncertainty will see this for the data points nearest neighbours. Thus, this singular data point suffers the potential of being surrounded by $N - 1$ other data points. The benefit this provides in fitting the model is thus extremely limited, and only slightly greater than if one data point had been chosen. A simple fix would be to simulate the model after 1 iteration, and select the next point from here. By doing this $N - 1$ times, a better solution may be found, although this may prove to be computationally very expensive.

2.3 Drug Data for Machine Learning

There are numerous data categories that can be used to represent a chemical in a suitable form for machine learning. Indeed, the field of chemoinformatics is dedicated to the pursuit of describing chemicals for computational models. Each of these methods have various strengths and weaknesses. Some are directly based upon the chemical structure whereas others are based upon physical properties. These can be combined to produce models with high predictive capabilities.

2.3.1 Physical Properties

A selection of physical properties from chemicals are known, from melting points to solubility. Many of these provide important aspects for consideration and allow human scientists to predict interactions, especially when determining new drugs. These data are often reported in tables within textbooks such as Perry's [] or provided through software [chembl ...]. Several of these data can be predicted through theoretical models, although the difficulty increases for larger molecules. For example, Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

2.3.2 Fingerprints

Another methodology is to develop a fingerprint: a unique code based on the chemical structure, either of the atomic arrangement, or by the electron cloud distribution. The latter of these is more fundamental to the activity of molecules but far harder to calculate. Indeed, for accurate representation of the latter, both atomic structure is needed *and* solutions for the Schrödinger equations corresponding to molecule in question.

According to Capecchi, Probst, and Reymond [CPR20], the most popular fingerprint in use are Morgan Fingerprints, a form of Extended Chemical Fingerprint (ECFP). ECFPs use a simple algorithm in order to generate a unique identifier, as described by Rogers and Hahn [RH10]:

1. **Initial Assignment:** Each atom has an integer assigned as an identifier.
2. **Iterative Updating:** Updating the identifier assigned to atoms based on adjacent atoms and structural duplications.
3. **Duplicate Removal:** Duplicate features are removed for hashing.

The iteration process involves each atom and adjacent atoms sharing numbers before in an array. A hash function is applied to this array and becomes the atoms new identifier.

Chapter 3

Methodology

3.1 Outline

The methodology presents a novel means of assessing different parametrised active learning methods on existing data sets, allowing for a robust answer into the use of active learning in drug rediscovery. Results can thus be given with a given belief. This approach has taken principles commonly used in machine learning and applied it to more traditional algorithmic methods.

Firstly, a collection of pre-existing data sets, X , are used. X is then split into two sub sets: X_{train} and X_{test} . Similarly to machine learning, the former of these subsets is used in fitting the parameters of the equation, and the latter is used to provide a result without the risk of data leakage into the training set. This is represented in [1]. Parallelisation is used to efficiently train the algorithms allowing the time for training to be $\sim \mathcal{O}(c)$.

Examining the smaller details, each algorithm is provided with the sets x_{known} , y_{known} , and x_{unknown} . Various algorithms are given these sets and allowed to generate a subset of x_{unknown} to be added into x_{known} alongside corresponding y_{known} . This can then repeat until a predefined stopping point is reached. Scores are reported using a weighted mean squared error [2] based upon y_{predict} for all x . This is similar to a standard machine learning methodology with a couple of differences. Firstly, no distinction is made between the training and testing set within a dataset contrary to standard practice. This is due to two reasons. Firstly, the datasets are not large enough for an accurate representation of the data within the testing set, and secondly, the scoring to each dataset is not used within the machine learning algorithms to fit parameters as is usually the case. All algorithms used rely upon a simple custom composite model to allow for flexibility and consistency.

In Section [1], it was discussed that there are various methodologies of representing chemicals and drugs. ... (if time)

3.2 Proof

In order to demonstrate the effectiveness, a few data sets are used instead, and the program is executed function by function. To start with, the underlying custom functions will be demonstrated, followed by the algorithms and then finally the training framework.

3.2.1 Custom Base Functions

Split

The split function allows for each dataset to be split into x_{known} , y_{known} , x_{unknown} , and y_{unknown} , as demonstrated in Figure 3.1. This is required as a fundamental step for the algorithmic testing. To demonstrate the validity of this function, ...

Repartition

Upon each iteration, the sets provided to the algorithms need to be repartitioned to allow for the continual operation of the algorithm. This consists of two parts: expanding the known sets and removing entries from the unknown sets.

Model

The machine learning model is the only custom class used. Here, a similar structure is used when compared with sci-kit's [] machine learning. To manage this, it has four methods: `__init__`, `fit`, `predict`, and `predict_error`. The last of these is not seen in all sci-kit's machine learning models and is reserved for those which can report a certainty of prediction. Here, this was achieved by taking a standard deviation of the models.

Validate

This is a simple method with greater potential than has been explored. By providing this as a separate method, a more computationally intensive validation model could be used without interference as parallelisation could be exploited. However, as it currently stands, it returns the weighted mean squared error using the standard method provided by [].

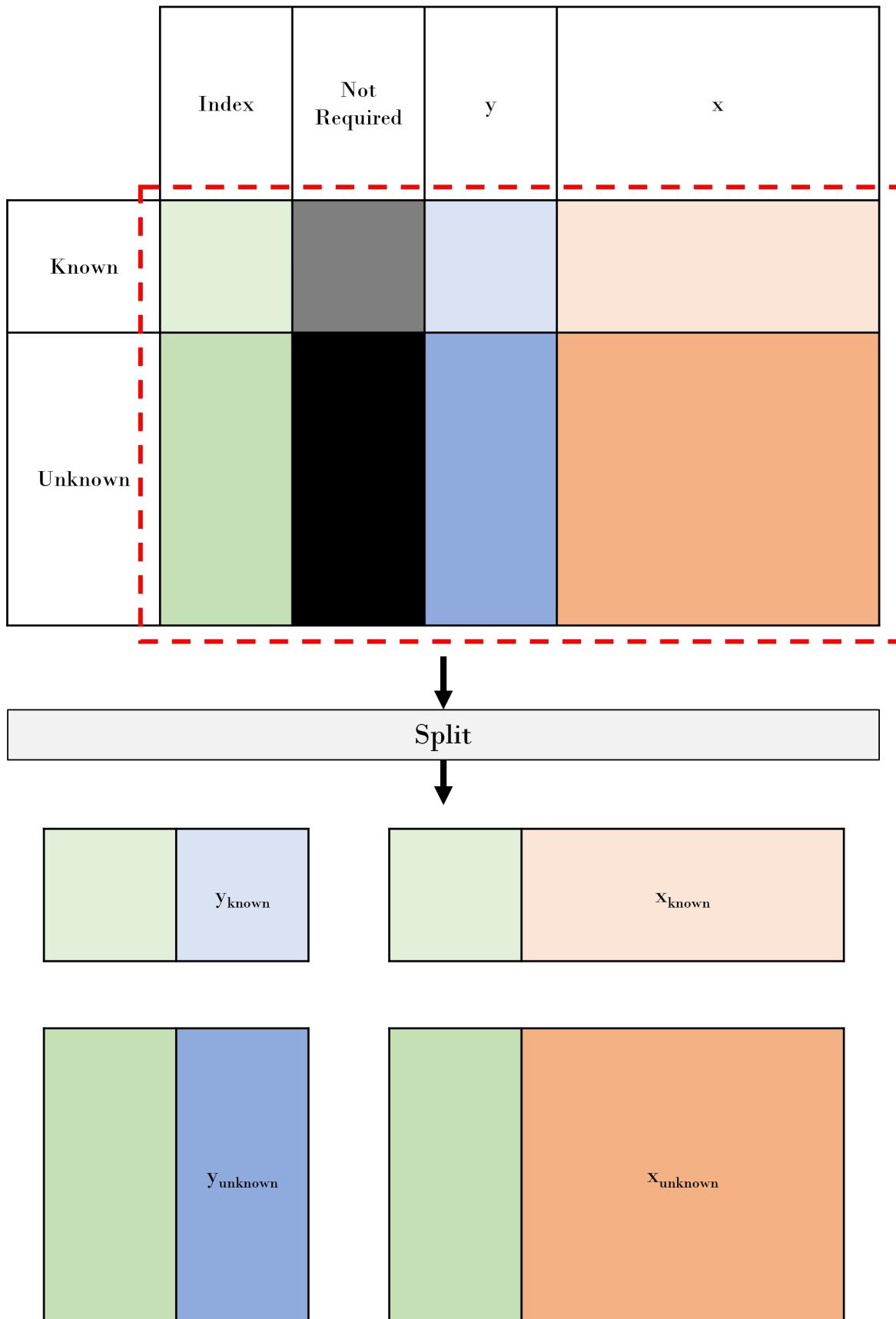


Fig. 3.1 Graphical representation of the split function. The red dashed boundary represents the input (additional colour coding has been performed to assist the reader in understanding the transposition of the base components).

3.2.2 Active Learning Algorithms

Dumb

The dumb algorithm, also referred to as random sampling or Monte Carlo sampling, refers to an algorithm that calls upon random samples to be tested. This represents the computationally least expensive approach, and is thus used as a baseline in comparing other algorithms. Since the datasets are shuffled prior to being used, the algorithm is extremely simple, as demonstrated in Algorithm 2.

Algorithm 2: Uncertainty Sampling Selection

Data: X_{unknown}

Result: X ordered according to priority for sampling

return X_{unknown}

Region of Disagreement

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

3.2.3 Training Framework

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Parallelisation

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Minimisation

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Chapter 4

Results

4.1 Overview

lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

4.2 Second Part

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

4.3 Special Case: COVID-19

Chapter 5

Discussion

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam,

luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Chapter 6

Conclusion

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam,

luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

References

- Capecci, Alice, Daniel Probst, and Jean-Louis Reymond (June 12, 2020). “One molecular fingerprint to rule them all: drugs, biomolecules, and the metabolome”. In: *Journal of Cheminformatics* 12.1, p. 43. ISSN: 1758-2946. DOI: 10.1186/s13321-020-00445-4. URL: <https://doi.org/10.1186/s13321-020-00445-4> (visited on 05/06/2022).
- Center for Drug Evaluation and Research (Apr. 25, 2022). “Coronavirus Treatment Acceleration Program (CTAP)”. In: FDA. Publisher: FDA. URL: <https://www.fda.gov/drugs-coronavirus-covid-19-drugs/coronavirus-treatment-acceleration-program-ctap> (visited on 05/05/2022).
- Pardo, Joe et al. (May 22, 2020). “The journey of remdesivir: from Ebola to COVID-19”. In: *Drugs in Context* 9, pp. 2020–4–14. ISSN: 1745-1981. DOI: 10.7573/dic.2020-4-14. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7250494/> (visited on 05/05/2022).
- Rogers, David and Mathew Hahn (Feb. 4, 2010). “Extended-Connectivity Fingerprints | Journal of Chemical Information and Modeling”. In: *Journal of Chemical Information and Modeling* 50.5. DOI: 10.1021/ci100050t. URL: <https://pubs.acs.org/doi/10.1021/ci100050t> (visited on 11/01/2021).
- Scikit Learn (2022). 2.3. *Clustering*. scikit-learn. URL: <https://scikit-learn.org/stable/modules/clustering.html> (visited on 05/05/2022).
- Settles, Burr (2009). *Active Learning Literature Survey*. Technical Report. Accepted: 2012-03-15T17:23:56Z. University of Wisconsin-Madison Department of Computer Sciences. URL: <https://minds.wisconsin.edu/handle/1793/60660> (visited on 11/01/2021).
- Settles, Burr and Mark Craven (Oct. 25, 2008). “An analysis of active learning strategies for sequence labeling tasks”. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. EMNLP ’08. USA: Association for Computational Linguistics, pp. 1070–1079. (Visited on 05/01/2022).