

Multi-output and Probabilistic Large Scale Benthic Habitat Mapping

JUSTIN TING

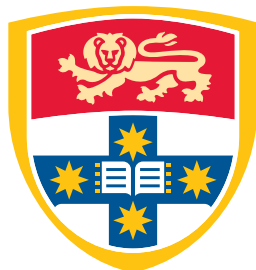
SID: 430203826

Supervisor: Dr. Simon O'Callaghan

This thesis is submitted in partial fulfillment of
the requirements for the degree of
Bachelor of Information Technology (Honours)

School of Information Technologies
The University of Sydney
Australia

29 September 2016



THE UNIVERSITY OF
SYDNEY

Student Plagiarism: Compliance Statement

I certify that:

I have read and understood the University of Sydney Student Plagiarism: Coursework Policy and Procedure;

I understand that failure to comply with the Student Plagiarism: Coursework Policy and Procedure can lead to the University commencing proceedings against me for potential student misconduct under Chapter 8 of the University of Sydney By-Law 1999 (as amended);

This Work is substantially my own, and to the extent that any part of this Work is not my own I have indicated that it is not my own by Acknowledging the Source of that part or those parts of the Work.

Name: Justin Ting

Signature:

Date:

Abstract

Being able to predict the state of benthic habitats based on limited information is crucial for environmental conservation, particularly as the impact of human activity on our oceans is greater than ever before. A considerable portion of work done in the area uses deterministic methods that strictly assign only one label to a given bathymetry data point, while more advanced models provide probabilistic results over all possible labels at any one point, also similarly only representing a single output. However, like the majority of real life classification problems (citation here perhaps), habitat mapping is intrinsically a multi-label problem for any data collected at a resolution low enough to be economically feasible to be performed at a large scale. In this paper, we explore advantages of having probabilistic class outputs as well as treating benthic habitat mapping as a multi-output problem, particularly when working with relatively low resolution bathymetry data, compared to the primary method of deterministic, single-output methods explored in existing literature.

Acknowledgements

The thanks go in here.

Contents

Student Plagiarism: Compliance Statement	ii
Abstract	iii
Acknowledgements	iv
List of Figures	vii
List of Tables	viii
Chapter 1 Introduction	1
1.1 Contribution	1
1.2 Motivation	2
1.3 Outline	2
Chapter 2 Literature Review	3
2.1 Overview	3
2.1.1 Habitat Characterisation	3
2.2 Habitat Classification	5
2.3 Map Creation	5
2.4 Non-Machine Learning Approaches	6
2.5 Machine Learning in Benthic Habitat Mapping	7
2.5.1 Deterministic Machine Learning Algorithms	7
2.5.2 Probabilistic Methods	9
Chapter 3 Experiments	13
3.1 Mathematical Background	13
3.1.1 Gaussian Processes	13
3.1.2 Dirichlet Multinomial Regression	13
3.2 Preprocessing	15
3.2.1 Downsampling the Data	15

3.2.2	Simplifying labels	17
3.2.3	Coordinates as features	19
3.2.4	Subsampling for Gaussian Process experiments	19
3.3	Illustrative Example	20
Chapter 4	Results	26
4.1	Deterministic Methods	27
4.2	Gaussian Process Classification	27
4.3	Dirichlet Multinomial Regression	28
Chapter 5	Evaluation and Discussion	29
5.0.1	Limitations	29
Chapter 6	Conclusion	30
6.1	Future Work	30
Bibliography		31

List of Figures

3.1	create Fixed-sized grids placed over training data	16
3.2	Dendrogram of training data	17
3.3	Distribution of labels in original dataset	18
3.4	Distribution of labels in multi-label outputs	18
3.5	Distribution of simplified labels in original dataset	18
3.6	Distribution of simplified labels in multi-label outputs	18
3.7	Full predictive map using Random Forests including coordinates as features	19
3.8	Full predictive map using Random Forests excluding coordinates as features	19
3.9	Plots of the three clusters, with labels taking on the argmax of each point	20
3.10	Legend/axes for the following histogram plots showing distribution of labels at each point	21
3.11	Label distribution of cluster A	21
3.12	Label distribution of cluster B	21
3.13	Label distribution of cluster C	21
3.14	DM Label distribution of label 0	23
3.15	DM Label distribution of label 1	23
3.16	OvR GP performance with variance for label 0	24
3.17	OvR GP performance with variance for label 1	24

List of Tables

CHAPTER 1

Introduction

Earth's oceans cover 70% of its surface, but only less than 10% of the Earth's oceans have been explored to date¹. There have been increasing efforts over the past few decades to more efficiently map out these unexplored areas to monitor marine ecosystems to be able to track the state of them over time for management, preservation, etc. purposes. The process used is called benthic habitat mapping, which is the process of generating predictive maps of different habitat types at the bottom of a body of water. Most studies looking to create benthic habitat maps share some basic key steps - acoustic data is used to estimate properties about the surface of the water, which are then mapped to, using machine learning algorithms, *in situ* data such as still images, videos, or samples of the area in question. It is the relationship which is inferred between the different data sets inferred using machine learning techniques that varies between studies. A considerable portion of such studies are shown to use deterministic methods to predict a label for any given coordinate such as Random Forests and Support Vector machines (SVMs), whilst more recent ones make use of more informative methods such as Gaussian Processes, providing a distribution over all possible labels given any data point.

1.1 Contribution

The main contribution of this thesis will be to explore how to use data where a single data point does not only have one label exclusively, but instead corresponds to a tally of each possible label. For example, a particular 5m x 5m area in the benthos may be an even mix of both sand and coral, but in previous literature, the data was simplified such that whichever label occurred more frequently regardless of how small the margin would be the single label assigned to that point. This results in a very coarse approximation even when using Gaussian Processes attempts to model the uncertainty/uncertainty with its predictions at each point (but ultimately only provides a single, final prediction). To alleviate this

¹Oceanservice.noaa.gov. (2016). How much of the ocean have we explored?. [online] Available at: <http://oceanservice.noaa.gov/facts/exploration.html>

and provide a richer set of information, we explore the use of Dirichlet Multinomials, which provides a distribution of each label that represents something entirely different. Whereas in a Gaussian Process, each label is assigned the probability of being the correct one, the output of a Dirichlet Multinomial Regressor provides the distribution of the frequency of labels in a particular space itself. See section **GP vs DM** for an illustrative example on how results would differ in practice between the two methods.

1.2 Motivation

The motivation behind assessing the effectiveness and advantages of such a method are that they inherently tie in with lower resolution data, particularly when a single images corresponds to a large enough area such that one would expect a mix of different labels. This is advantageous because we want to be able to re-sample data from any given site periodically (for example, every 3-4 years) whilst being economically efficient. This naturally lends to lower resolution data, meaning that summarising large areas to a single label would theoretically be throwing away a majority of the information contained in bathymetry and image data.

1.3 Outline

outline

Literature Review

2.1 Overview

The process of benthic habitat mapping involves three key steps that the large majority of all studies in the area go through.¹ In this section, we will give a brief overview of each of these steps, along with common procedures used in them across studies in this area.

- (1) **Habitat Characterisation** - extracting properties of the environment such as rugosity (roughness), aspect (direction of slope), depth
- (2) **Habitat Classification** - grouping the raw information about the environment into categories, such as sand, granite, etc.
- (3) **Habitat Mapping** - using classifications with the larger scale bathymetry data to extrapolate habitat maps

2.1.1 Habitat Characterisation

not just resolution of data but modality. images vs bathymetry If we were able to collect high resolution data for the entire ocean's benthos - the job of creating benthic habitats for any given area would be (relatively) trivial. As this is prohibitively expensive, we instead collect large amounts of low resolution data, and small samples of high resolution data (between which we model a relationship). This subsection provides a brief summary of data collected and methods used to do so.

Remote-sensing data. Due to the cost of sea expeditions, it is economically infeasible to have marine vehicles (autonomous or otherwise) explore the entire ocean floor to confirm the ecological properties of all of Earth's benthos. However, we do need to collect sufficiently detailed data of large areas at a time,

¹Ozcoasts.gov.au. (2016). Benthic habitat mapping: Mapping Overview. [online] Available at: http://www.ozcoasts.gov.au/geom_geol/toolkit/mapoverview.jsp

particularly those of being mapped, and for this, remote-sensing data is used. These usually come in the form of acoustic backscatter data that involves the firing of sound waves towards the benthos, whereby their frequency and strength upon returning is used to deduce the depth of a particular material, as well as the density of said material (from which a guess at the actual substance can be made - e.g. sand, mud, etc.).

Multibeam echosounders (MBES) are becoming a more frequently used method of collecting acoustic backscatter data (Calvert, Strong, McGonigle, and Quinn, 2015) despite older methods involving single beam echo sounders (SBES) being cheaper and easier to segment. This stems from the fact that the reduced cost comes at the expense of (potentially) accuracy, as well as lower resolution data. This is due to SBES' beam angle, i.e. the angle formed by the 2D flattening of the 'cone' shape of the emitted beams, ranging from 15-25°, whereas MBES' is 0.5-3°, depending on the particular system (Brown, Smith, and Lawton, 2011). The difference in angle means that data returned via SBES devices are more 'coarse', representing less accuracy and granularity, whereas that of MBES is more detailed and can present more information. However, there is overhead associated with use of MBES, in that the considerably decreased angles means much more 'overlapping' data, adding complexity to the segmentation process.

Truthing Data. explain redundancy here, unclear what it's referring to - why is there redundancy?

The most common methods to be able to obtain a sufficiently large truthing data set (but still trivially small compared to the area covered by remote-sensing data) are videos or images - though the former still requires post-processing to extract the needed images. The advantage that can be provided here, however, is the redundancy in data points (Rattray, Ierodiaconou, J. Monk, and Kennedy, 2014) - but there is extra cost in time required to convert videos into the needed images (pre-processing before feeding into algorithms for habitat mapping), an area that is in itself worth of research within the field. (Lucieera, Hilla, Barretta, and Nichol, 2013)

Other data. (why is water column correction important when correlating images with seagrass standing crop?) Other data that is less common, but also used to map habitats, is patterns in the water movement (such as tidal currents, wave action) (Brown, Smith, and Lawton, 2011) in the column of water above the area of benthos being mapped - a feature that provided useful input in arriving at an accurate benthic habitat map (in addition to sediment analysis). (Snelgrove, 1994) Other sources such

as UNESCO have also verified the importance and significance of using water column correction techniques to obtain more accurate habitat maps, particularly when correlating images with seagrass standing crop.²

2.2 Habitat Classification

more in-depth focus here, what kind of supervised/unsupervised ML algorithms are used for classification? Almost all studies use *in situ* 'truthing' data to complement the acoustic data to be able to build a model between the acoustic data and truthing data (creation of these models are explained in following sections). However, we need to know the labels of this data considering that the final goal is to create a habitat map, where any one habitual zone is given its prospective label - to do this, we also need to label the clusters of truthing data. These categories may be, for example, 'bedrock covered by discontinuous seagrass cover', 'Maerl interspersed with sand and gravel', 'superficially coarse sand to fine gravel covered by dense patches of seagrass', etc. (Micallef et al., 2012). The two overarching ways to perform this classification are in the form of supervised and unsupervised algorithms.

Studies have used both supervised and unsupervised methods in clustering the initial data for the training step. Often, there may be large amounts of visual data, beyond that which any human or even team can reasonably, manually cluster - and as such, unsupervised algorithms are first used to create these clusters, after which an expert may be brought in to verify/simplify (or otherwise) the resulting clusters. (Steinberg, Friedman, Pizarro, Williams, and S.B., 2011)

2.3 Map Creation

The final step is map creation, which many papers related to benthic habitat mapping focus on - and also where the most variation occurs in terms of the method used. The various approaches used can be categorised into two broad categories. The first is a top down approach whereby the classification of the habitat characterisation data is validated (or otherwise) with the truthing data, and the second is a bottom up approach where the characterisation data is similarly clustered into classes, but not to directly represent a particular habitat - instead, the aim is to find a relationship between the acoustic data clusters and the truthing data clusters which we can model. Using this model, we can then extrapolate

²Unesco.org. (2016). Water column correction techniques. [online] Available at: <http://www.unesco.org/csi/pub/source/rs10.htm>

the acoustic data which doesn't have corresponding truthing data to create the habitat map. (Ahsan et al., 2011) We will explore this aspect more when looking at how the mapping process has evolved over time and the improvements that it has brought about.

2.4 Non-Machine Learning Approaches

this section should be part of the previous one (Map Creation) While the majority of modern papers in benthic habitat mapping employ machine learning techniques for map creation, this doesn't exclude those that do not from providing useful information and insight. An important study was undertaken in 2001 that employs relatively basic statistical analysis, employing different forms of variance as its main tool of analysis. (Kostylev et al., 2001), at the time (and in fact, even now) integrated more sources of data together than most other studies undertaken - multibeam bathymetric data, geoscientific data, seafloor photographs, habitat complexity, and relative current strength. Rather than drawing broad conclusions about the effectiveness of a collection of tools in creating habitat maps, deeper analysis is done on subsets of the data to attempt to clarify some of the complexities and intrinsic properties of benthic habitats and ecosystems themselves. Although little is done to address and verify accuracy of the actual results/map in this paper, it provides value through the analysis of variance and covariance performed on and between different benthic/marine properties, establishing relationships typically taken for granted or ignored. For example, it is established that while sediment type contributed heavily to a higher taxonomic group count, there was little relationship between sediment type and depth. However, this only indicates that there is no linear relationship between the two, and doesn't necessarily preclude a non-linear relationship between them, perhaps with the inclusion of other parameters as well. In particular, Kostylev establishes that gravel substrates are more abundant with varying taxonomic groups than their sand counterparts.

Certain organisations, government bodies/etc. will also provide guidelines outlining the classification process. For example, the European Nature Information System website and the Australian Government's 'Interim Marine and Coastal Regionalisation for Australia'³ both provide classification schemes for people creating habitat maps or other similar efforts.

Such findings provide useful insights for future studies that will allow a better assessment of data, or to be able to apply initial assumptions in obtaining better results. However, constantly seeking a deeper

³Unesco.org. (2016). Water column correction techniques. [online] Available at: <http://www.unesco.org/csi/pub/source/rs10.htm>

understanding through a proportionally increasing amount of sampling creeps towards 'exploring' the entire Earth's oceans manually. To obtain economically feasible yet reliable predictions from the limited data that we have, we need to employ machine learning techniques to fully utilise the information that we gather.

2.5 Machine Learning in Benthic Habitat Mapping

As benthic habitat mapping covers a diverse range of disciplines, namely "marine biology, ecology, geology, hydrography, oceanography and geophysics" (Brown et al., 2011), in addition to statistics and machine learning, it is logical that it would take considerable effort and vast resources to give each relevant discipline an equal, and large amount of attention within any single study. Thus, different papers can rely on collective findings of others to launch their own research and look further into particular lines of inquiry, start new ones altogether, or evaluate effectiveness of methods used in the field/etc. Within benthic habitat mapping, one prevalent line of inquiry is how to create more accurate, higher quality maps by employing machine learning techniques. For the remainder of this review, we will be revisiting common machine learning techniques and their application in the various stages of benthic habitat mapping along with the benefits they provide.

2.5.1 Deterministic Machine Learning Algorithms

In this section, we will review some machine algorithms that can be used in benthic habitat mapping processes - whether that be in the initial clustering stages of (ideally) independently gathered datasets such as acoustic backscatter data and collections of high resolution images, or the actual classification of 'new' (or testing) data in determining their predicted habitat classes.

Multinomial Logistic Regression. Multiple Logistical Regression is one of the more basic machine learning algorithms that can be used to predict habitat classes, and falls under the 'supervised learning' category as we have the 'output' for the feature vector in the initial data. Regression, broadly, involves the estimation of relationships between variables, and logistic regression involves the prediction of likelihood of class membership given a number of variables (that are assumed to have low collinearity). This only applies to domains with two classes, however - to use this technique for classification where we have an unbounded (though usually still relatively low) number of classes, we need to use multinomial logistic regression, which is able to account for more than two distinct, unordered (i.e., sand vs. mud has

no relative ordering) classes, where class membership is predicted using maximum likelihood estimation (MLE), similarly to logistic regression. However, the difference is that whereas logistic regression only requiring a single logit function as its nominal variable is dichotomous, multinomial logistic regression requires comparison between $k - 1$ (where k is the number of possible dependent variables) logit functions.

Even though Caruana and Niculescu-Mizil (2006) show that logistic regression methods achieve on average worse results than most other approaches available, it recognises that in certain cases the models that perform most poorly on average still display exceptional performance, and as such, this method is still worth exploration and experimentation. In particular, Belanger et al. (2012) used multinomial logistic regression across temperature, salinity, and productivity to correctly predict class membership by a margin of 23-84% more than by pure chance. This is equivalent to an improvement of 1-2x compared to a random guess, which taken at face value would suggest that logistic regression is an undesirable choice of algorithm for this problem domain.

Random Forests. In contrast to logistic regression, random forests were shown in Caruana and Niculescu-Mizil (2006) to be state of the art, only just falling short of boosted decision trees after calibration. Random forests are an ensemble method, meaning that it uses a collection of estimators, before aggregating their results to obtain some sort of average. The aim of this is to minimise the variance and hence error that any single one of these estimators would otherwise result in.

From the initial dataset, some number B is chosen which represents the *number* of trees to build (as a part of our random forest), after which, B random, unique subsamples of the full dataset are taken. Within each decision tree in our random forest, some constant number m of features is taken at each node of the tree, such that the split at each node only takes into account the m randomly chosen features. Each of the decision trees in our forest will hence have a 'result' (that may be a class or some continuous value). Typically, the final decision of the random forest will be made by a vote count for classification, and an average of each decision tree's result in regression problems.

As random forests are a method that is low in complexity but provides very good results on average, we can see that it is used in quite a few studies (Lucieera et al. (2013), Seiler et al. (2012), Hasan et al. (2014)), where the random forest classifier provided the best results over other methods relating to at least a significant subset of the explored data. However, Lucieera et al. (2013) found that while random forest classifiers were most able to classify substratum and rugosity, K-nearest neighbour classifiers most

accurately classified sponge structure classes, pointing to the need to do a more systematic comparison of different methods in benthic habitat mapping. A further advantage to using random forests as pointed out in (Hasan et al., 2014) is that it can provide insight into which features were more important than others, which can aid future studies to be more successful and efficient by focusing more efforts towards collecting the most influential data. The success met with using random forests make it a good benchmark to compare against for future work that aim to develop methods to create more accurate benthic habitat maps than has been done before.

Multi-class Support Vector Machines. Although support vector machines fell outside the top three overall supervised learning algorithms in terms of performance (accuracy), their non-parametricity could potentially be of benefit given that our knowledge of the complex relationships between elements of benthic habitats are limited. Moreover, despite SVMs being rarely used anywhere in the field, they are "acknowledged to be very competitive discriminative classifiers in machine learning literature" (Ahsan et al., 2011).

However, SVMs in their base form only support classification into two classes, requiring modification to the original algorithm to support more - an active area of research that has not found any single 'best' way to perform this algorithm extension yet. To do so, there are two, basic main approaches available, one being a **one vs. all** approach, and the other being a combination of all **one vs. one** approach. We can hence quite clearly see that given C possible classes, the first would require C separate classifiers, whereas the latter would require $\frac{C(C-1)}{2}$ classifiers (Murphy, 2012). Multi-class SVMs were used in Ahsan et al. (2011) for illustrative purposes, using the one vs. one approach as per their use of LibSVM (Chang et al., 2011), outperforming classification trees on certain datasets, and in limited cases (but not overall), Gaussian Mixture Models as well. Again, the mixed results would suggest that under certain conditions (such as size of the dataset and various properties of the data itself, some of which are not known before testing), use of a multi-class SVM could provide a useful benchmark to some extent.

2.5.2 Probabilistic Methods

The classifications being made regarding benthic habitats naturally involve uncertainty, as we are still learning the relationship between different characteristics of benthos with the varying communities of

fauna and flora that reside there. Whilst guessing the most likely class for a particular domain deterministically has its practical applications, it is arguably more *natural* to represent the uncertainty (Rasmussen and Williams, 2006). As our understanding of marine environments is still quite weak (of the United Nations, 2004), it is debatable whether deterministic results are always appropriate when being used to make high level management decisions relating to marine environments. While deterministic methods will create a model that attempts to explicitly account for all variables, probabilistic models deal with joint distributions over all the variables. As we need to better understand "the complexities of coastal system functioning rather than simplifying and scaling down the system into smaller components" (Diaz et al., 2004), this feature can be especially valuable seeing as there is simply not enough 'expert knowledge' to adequately, explicitly model the relationship across a range of variables.

Illustrative Example. not an illustrative example - give some actual figures/graph of crossover point

A simple example of this can be seen when comparing the deterministic approach of a logistic regression classifier, with the probabilistic Naive Bayes classifier. Starting from no data, up until a certain threshold, a Naive Bayes (NB) classifier will actually provide a more accurate classification as it approaches its comparatively higher asymptotic quicker, after which point, once there is sufficient data, the logistic regressor will provide the better results (Ng and Jordan, 2002). In this simplified example, an analogy can be drawn where the data used up to the threshold when the NB classifier performs better represents a lack of knowledge about the data causing the logistic regressor to underperform, whereas the continued addition of data represents more understanding (more data points) of the domain, allowing logistic regression to then outperform the NB classifier.

Gaussian Mixture Models. Gaussian mixture models (GMMs) are parametric models that "model the distribution of data as a set of clusters, where each cluster is a multivariate Gaussian" (Ahsan et al., 2011). In this particular paper, GMMs are compared with classification trees, which it is found to perform better than in most cases, but were also predicted classes from unseen data with higher certainty than discriminative methods. This is because of its generative nature that accounts for the distribution of bathymetric features, allowing it to model the joint distribution of the classes as well as features. Moreover, each (function = distribution) Gaussian function within the model has its own mean and covariance matrix, which also contributes to its powerful modeling ability. However, the use of GMM may have been hindered by the dimensionality of the data - while only five properties were measured, each was calculated for a varying number of scales for the input vector, meaning the 'features' were at least some multiple of five. As this exceeds the recommended six dimensions for use with GMM, application to a

very large dataset may be beyond reasonable computational ability.⁴ To avoid this, the feature vector may have to be truncated to contain the bathymetric properties for only one particular scale at a time.

Using Gaussian Processes. A recent study used probabilistic methods to develop a mapping between the clustered acoustic data to continuous cluster probabilities, as opposed to discrete cluster labels, thus representing the certainty of the results obtained. Using Gaussian Processes which do not inherently support classification, Bender et al. (2012) extended the probabilistic least squares classifier to retain the information regarding certainty of class membership that exists during the classification process, rather than discarding it in the traditional method. By evaluating the probabilistic results of PTLSC by comparing its results with the actual cluster probabilities obtained in the classification of the images via an unsupervised variational Dirichlet process model, it was shown that the PTLSC method performed better than a PLSC trained directly on the discrete cluster labels in terms of accuracy, mean squared error, and mean variance as well. This demonstrates that while both PTLSC and PLSC err in their predictions when dealing with the transition different boundaries, by maintaining probabilistic information in the PTLSC, it is able to make slightly better judgements in such cases.

Gaussian Processes and large datasets. However, Gaussian processes involve a matrix inversion process that requires an $O(n^3)$ operation which does not scale well with large datasets. To overcome this whilst reaping the benefits of Gaussian processes, Bender et al. (2012) extracted subsets of the original dataset on which to perform analysis - a small, randomly chosen portion from three Gaussians, of the initial millions of observations. While this has still provided a high accuracy for all methods tested, there is likely information to be gained by being able to use a considerably larger portion of the dataset. To do this, a method would be required to generate sparse covariance matrices through approximations (Bickel and Levina, 2008), or use of functions that guarantee sparseness as a property (Melkumyan and Ramos, 2009) - something that can be explored in future work. To illustrate how the obstacle can be overcome, the latter paper describes a method whereby, rather than inverting the covariance matrix in its raw form, a threshold is calculated at which point, rather than observing the normal 'tapering' off of covariance values, they are simply set to zero beyond that point. This will result in a significant portion of the covariance matrix being populated with 0s, at which point inversion of the sparse matrix can be performed for which there are known efficient methods. However, there have been more ways of sparse approximation GPs that other studies have explored.

⁴Nickgillian.com. (2016). GMM Classifier. NickGillianWiki. [online] Available at: <http://www.nickgillian.com/wiki/pmwiki.php/GRT/GMMClassifier>

Sparse Approximation Gaussian Processes. m isn't clarified here, plus $n < m$ is incorrect, should be $m \ll n$

lots of things here have become irrelevant (talks about 'optimal' methods that aren't implemented/included in experiments - limit to relevant ones, i.e. the GP ensemble methods

It is of importance that a number of methods of dealing with sparse approximation of GPs are taken into account if the aim is to deal with large GPs in the inversion step. (Quinonero-Candela and Rasmussen, 2005) explores exactly this, immediately discounting the "subset of data" (SoD) method as being non-competitive due to it not being able to represent the original data to a reasonably accurate enough extent, though we have seen that this was the approach taken in (Bender et al., 2012). As all the different methods (bar SoD) have a complexity of $O(nm^2)$ where n is the size of the data, and $n < m$, the authors notably point out that no gross approximations should be made as more competent methods are computationally equivalent, and as such point towards their notes on future work to outperform the existing state of the art. As such, we would wish to explore combining the **Partially Independent Training Conditional approximation** with "the most powerful selection method for the inducing inputs."

Experiments

To identify whether the Dirichlet Multinomial Regression method proposed can provide richer and more valuable information than single-output or deterministic methods can alone, we ran experiments on the data obtained from the ACFR's Sirius AUV and Schmidt's Falkor. The main machine learning algorithms' performance which we tested were Gaussian Process Classification, and Dirichlet Multinomial Regression. In this section, the experiments were designed to display the benefits of a Gaussian Process Classifier's probabilistic output, as well as the label distributions of a Dirichlet Multinomial Regressor.

3.1 Mathematical Background

(this needs/should be in its own chapter)

3.1.1 Gaussian Processes

GP equations, description

GP approximation methods

3.1.2 Dirichlet Multinomial Regression

Dirichlet multinomial regression, as the name suggests, combines dirichlet and multinomial distributions to achieve the combined model. In particular, we are interested in modeling a distribution over category counts, as there exists relationship in our data such that every bathymetry point corresponds to a certain count of each possible label in the relevant area of benthos. **explain why we should first revisit dirichlet, multinomial distributions separately before looking at dirichlet multinomial regression**

3.1.2.1 Multinomial Distribution

equations, description

3.1.2.2 Dirichlet Distribution

descriptions

$\theta \sim \text{Dir}(\alpha)$, dirichlet distributed random variable

$$p(\theta) = \frac{1}{\beta(\alpha)} \prod_{i=1}^n \theta_i^{\alpha_i-1} I(\theta \in S) \text{ density function, } I \text{ is indicator function}$$

$\theta = (\theta_1, \dots, \theta_n)$, $\alpha = (\alpha_1, \dots, \alpha_n)$, $\alpha_i > 0$ theta - n-dimensional vectors, alpha - parameters for distribution

$S = \{x \in R^n : x_i \geq 0, \sum x_i = 1\}$ S is probability simplex, the set of pmfs on numbers 1 through n

$$\frac{1}{\beta(\alpha)} = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_n)}, \alpha_0 = \sum_{i=1}^n \alpha_i \text{ generalised beta function}$$

3.1.2.3 Dirichlet Multinomial Regression

descriptions

$$DM(C|\alpha) = \frac{M!}{\prod_k C_k!} \frac{\Gamma(\sum_k \alpha_k)}{\Gamma(\sum_k c_k + \alpha_k)} \prod_{k=1}^K \frac{\Gamma(C_k + \alpha_k)}{\Gamma(\alpha_k)}$$

$$M = \sum_k c_k$$

For the regressor, the two activation functions that were considered were exponential and softmax, where the former often provided better mapping predictions, but the latter is preferable in the general case due to its better numerical stability [include graphs of exponential and softmax here](#).

$$\alpha_k = \exp\{x^T w_k\}$$

$$\alpha_k = \text{softmax}\{x^T w_k\}$$

The weights w here are in fact a matrix of weights with dimensions $(K \times D)$, where K is the number of possible labels across the dataset, and D is the dimensionality of the dataset. Multiplying the dirichlet multinomial prior by the likelihood then gives the equation over which to optimise to predict the normalised label counts at any given point.

This gives the joint-log-likelihood over both the dirichlet and multinomial distributions:

$$\begin{aligned}
& \sum_{n=1}^N [\log(M_k) - \sum_k \log(c_k!) + \log \Gamma(\sum_k \alpha_k(x_n)) - \log \Gamma(\sum_k c_{nk} + \alpha_k(x_n))] \\
& + \sum_{n=1}^N \sum_{k=1}^K [\log \Gamma(c_k + \alpha(x)) - \log \Gamma(\alpha_k(x_n))] \\
& + \sum_{k=1}^K [-\frac{\phi}{2} \log(2\pi\phi) - \frac{1}{2} w_k^T \phi \mathbb{I} w_k] \quad (3.1)
\end{aligned}$$

To optimise this equation, the partial derivative of the above over the weights w are considered:

$$\begin{aligned}
\partial \frac{\log p(c, x)}{\partial w_k} &= \sum_{n=1}^N x_n \alpha_k(x_n) [\psi(\sum_l \alpha_l(x_n)) - \psi(\sum_k c_{nk} + \alpha_k(x_n))] \\
& + \sum_{n=1}^N x_n \alpha_k(x_n) [\psi(c_{nk} + \alpha_k(x_n)) - \psi(\alpha_k(x_n))] - \frac{1}{\phi} w_k \quad (3.2)
\end{aligned}$$

explain all the symbols here

3.2 Preprocessing

3.2.1 Downsampling the Data

As the purpose of using Dirichlet Multinomial Regression was to be able to model the distribution of habitat label occurrences over an area, we downsampled the combined 2011+2015 dataset which was at a significantly higher resolution than the 2009 dataset. Two methods of downsampling in particular were tested. The first coarser approach involved simply taking the space in which the data was collected and placing grids of fixed size over them as in fig. 3.1, binning all points falling within each grid into a single datapoint. Each of these data points contained multiple points from the original dataset with their own counts for each of the possible labels, so the downsampled points simply took the sum of all the label counts in each fixed grid.

The second summed label counts in the same way, but clusters were instead formed by first calculating the full dendrogram on the 16502 entries in the training data, and forming groups such that none had more than 5 of the original points within them, and the sub-clusters (at each level of the dendrogram)

were no more than a 21 metres away from one another. As can be seen in fig. 3.2, the gradual merging into the single supercluster was quite consistent, indicating the original datapoints were mostly evenly distributed.

For a fair comparison between Gaussian process classification and dirichlet multinomial regression, the downsampled data was used to train the GPs as well - although this seems like an unnecessary handicap to the GP, it is more appropriate considering that one of the aims here is to demonstrate what sort of information can be gained from a DM vs. a GP, given the same *raw* data.

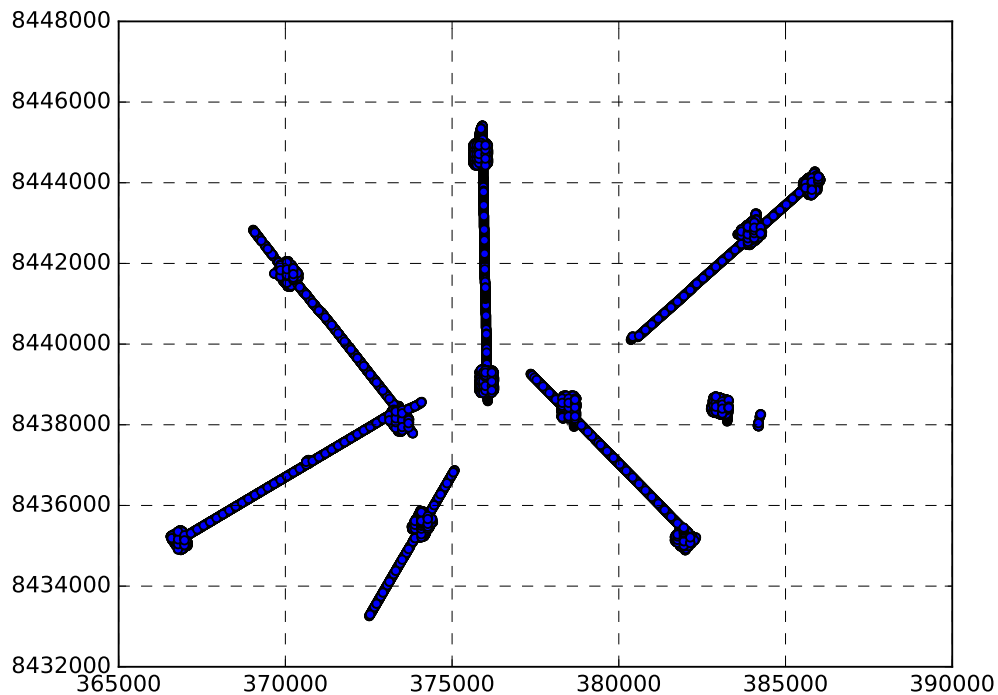


FIGURE 3.1: **create** Fixed-sized grids placed over training data

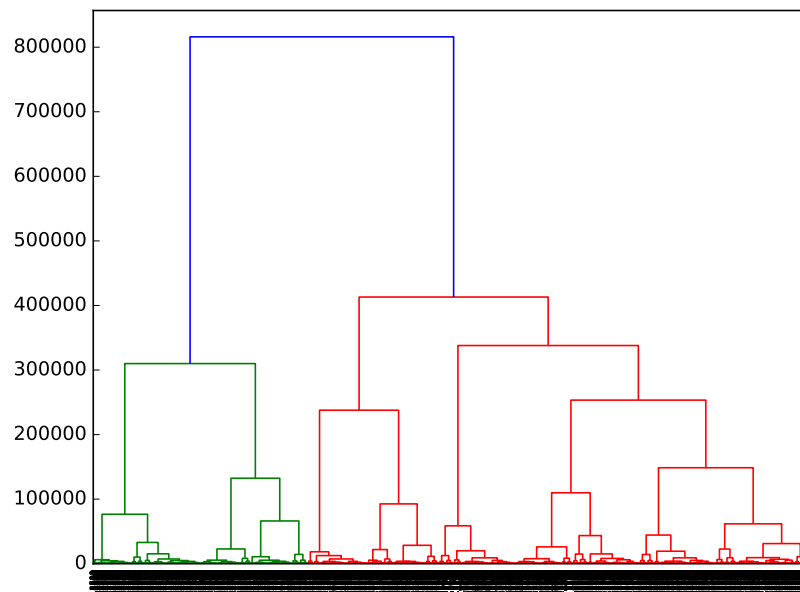


FIGURE 3.2: Dendrogram of training data

3.2.2 Simplifying labels

Another step that was considered during experiments was the aggregation of habitat labels. The original training data contained 24 separate labels determined through an automated clustering procedure using Dirichlet Processes. Because of the uneven distribution of these labels (generate these images 3.5 and 3.6), with the occurrence of some too insignificant for any machine learning algorithms to pick up, they were simplified in collaboration with ecological experts, who manually identified which of the 24 labels were in fact of the same class - for example, 5 separate classes of coral may have been indistinguishable to the average person, and were hence grouped into a single label. This allowed the near-non-occurring labels to be grouped together with more commonly occurring ones, whilst also allowing a different level of granularity in training models/forming predictions that could be used if only a rough approximation of an area's benthic map were required.

label mappings - give the labels for the simplified classes, e.g. coral, etc.

simplified	original
0	1, 2, 18, 20, 21, 23, 24
1	3, 5, 10, 16, 17, 19, 22
2	13, 14, 15
3	4, 6, 7, 8, 9, 11, 12

put some images here from both original and simplified classes. don't use squidle's downloader, visit server directory and d/l from there directly, 100x+ faster

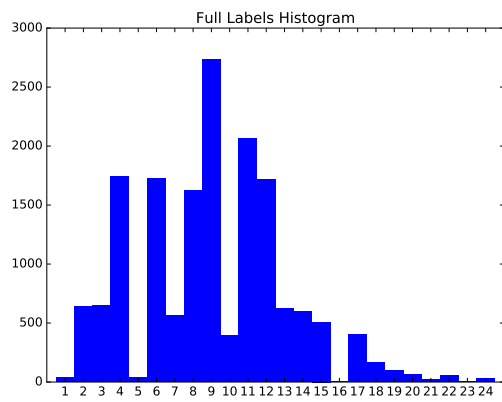


FIGURE 3.3: Distribution of labels in original dataset

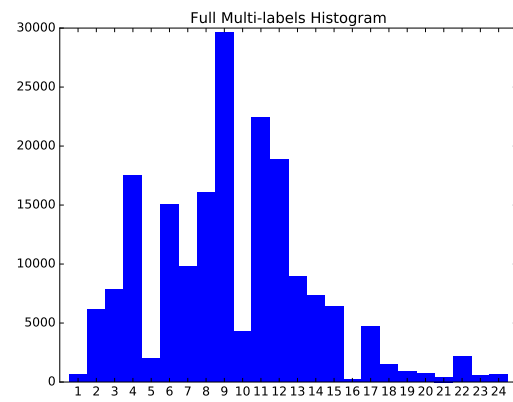


FIGURE 3.4: Distribution of labels in multi-label outputs

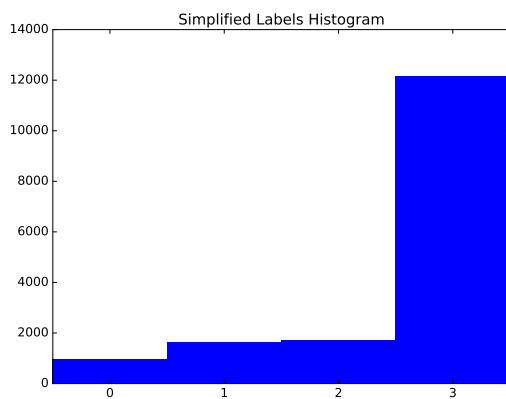


FIGURE 3.5: Distribution of simplified labels in original dataset

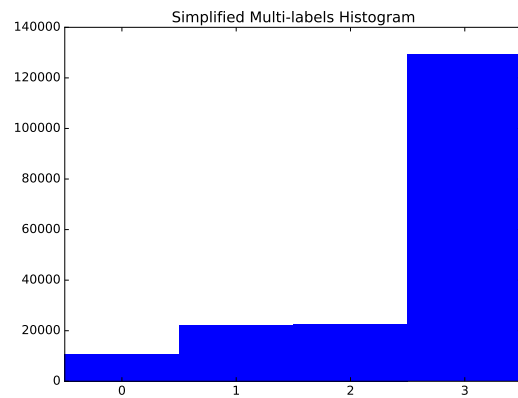


FIGURE 3.6: Distribution of simplified labels in multi-label outputs

3.2.3 Coordinates as features

Due to the abundant bathymetry data that was available in the form of depth, rugosity and aspect at each available data point, there was reason not to include the coordinates themselves in the feature space. Whilst it does make sense that in a natural environment, areas that were spatially near to one another would also have similar properties, this should not be relied upon, and other intrinsic properties should be the basis upon which predictions are made. Forming predictions on the full query dataset using a random forest supports this notion quite strongly - whilst 10-fold cross validation using the coordinates as features had a notably higher F-score of 0.61 compared to 0.40 without, the unnaturally straight split between the left and right segments over a 12km region suggests that the predictive map is flawed. (argument here alone is weak. for simplified labels using coords is still much better by a similar margin, do some reading to back this up properly)

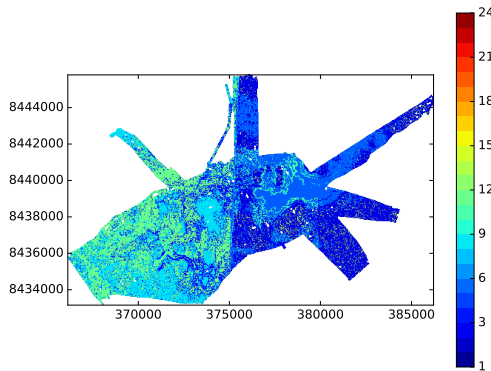


FIGURE 3.7: Full predictive map using Random Forests including coordinates as features

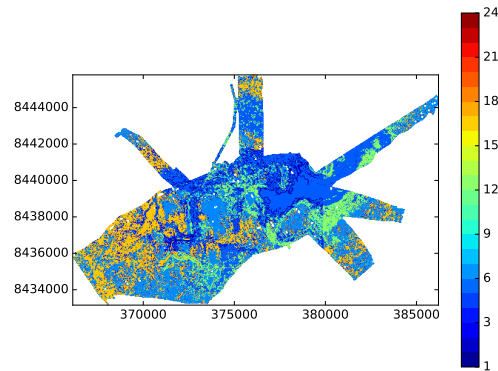


FIGURE 3.8: Full predictive map using Random Forests excluding coordinates as features

3.2.4 Subsampling for Gaussian Process experiments

Due to the $O(n^3)$ complexity of training a Gaussian Process Classifier, using all 16502 points was infeasible, so it was necessary to use only a subsample of the training data. As can be seen in the above histograms (reference the figure instead. may need to combine them into one), the distribution of classes in both the simplified and non-simplified versions was very uneven. As a result of this skew, randomly sampling the the training data to fit our GP classifier against resulted in worse results than sampling an equal *number* of points for each class. To obtain a reasonably well-performing set of 1000 points (the

number chosen to obtain a balance between performance and time required), 10-fold cross validation was performed on random sets of 1000 with each class sampled equally, and the best set chosen after 200 runs of random subsampling.

3.3 Illustrative Example

The differences between a Gaussian Process that provides the probability distribution of possible labels compared to the Dirichlet Multinomial Regressor that provides the distribution of actual labels at a point, are highlighted in the illustrative example below. Note that three clusters were synthesised, with clusters A, B containing 0.7 : 0.3 and 0.3 : 0.7 average ratios in label mix per point respectively, while cluster C contained an even 0.5 : 0.5 average split, where cluster had 100 points. The colours on the overall plot are only representative of the **most** common label at each point - the actual distributions at each point are shown in the graphs following it.

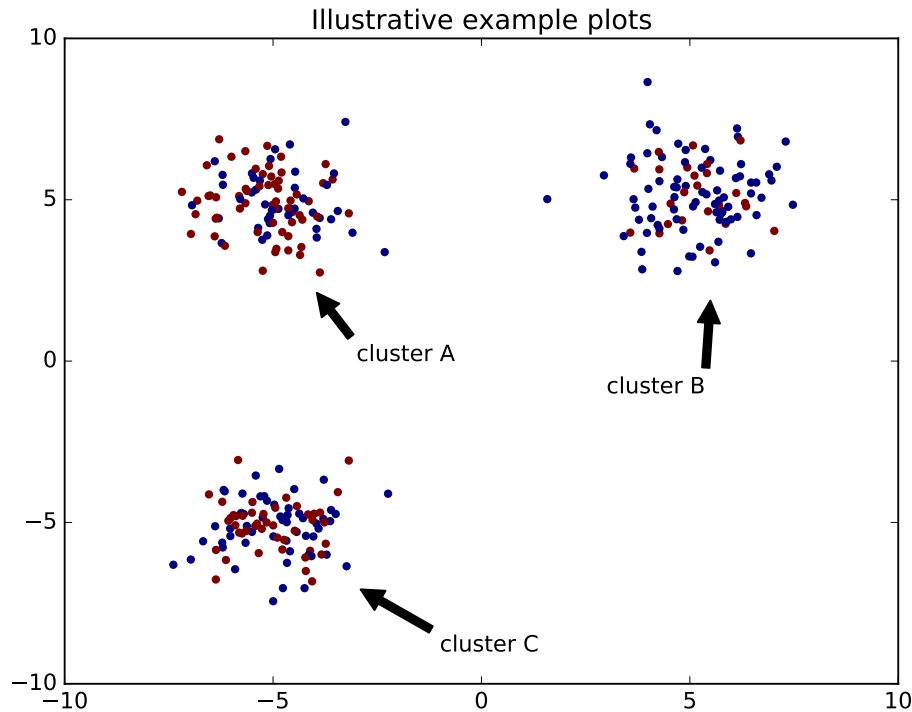


FIGURE 3.9: Plots of the three clusters, with labels taking on the argmax of each point

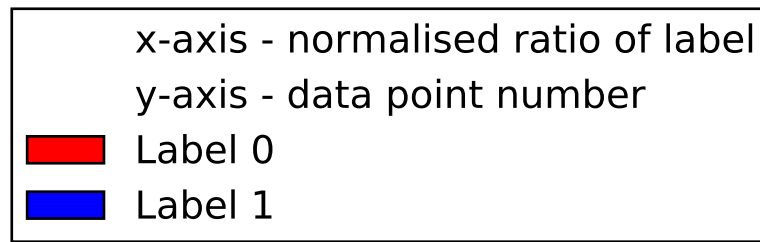


FIGURE 3.10: Legend/axes for the following histogram plots showing distribution of labels at each point

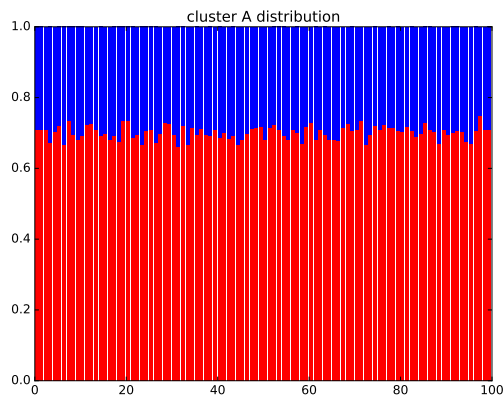


FIGURE 3.11: Label distribution of cluster A

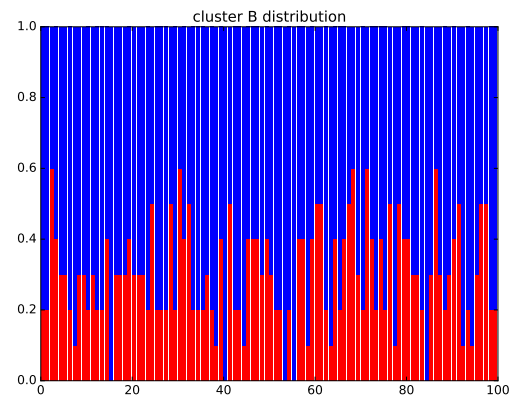


FIGURE 3.12: Label distribution of cluster B

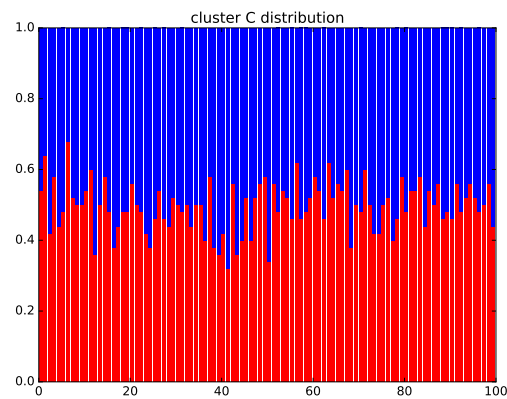


FIGURE 3.13: Label distribution of cluster C

In this example, the GP and DM models were each trained on half of each cluster, and made to predict the other half. However, as a standard GPC can only have single label inputs and outputs, a approximation/simplification was made for the purpose of calculating average error, whereby the label was simply taken to be the most frequently occurring label at any given point. While this is a reasonable simplification for clusters A, B as the dominant label has majority share, this is not the case for C, as the split between the two labels per point in the cluster is exactly even. In an initial attempt to counter this, multi-task GPs were considered as a means of making a *fairer* comparison between a GP and DM, but the idea was ultimately discarded as it was not fit for purpose, one of the primary issues being that the model does not inherently restrict the outputs of a given datapoint to sum to 1, instead being at the mercy of the parameters of the GP. The results and plots for this example are below, and figures displayed were taken from an average of 20 runs.

	Dirichlet Multinomial Regression RMSE*	Gaussian Process Classifier (argmax) RMSE
Original data	0.070179271314358999	0.2683333333333337
Quadratic-space projection	0.065630111843395234	0.4343333333333335
Cubic-space projection	0.29019235800882354	0.43725490196076466

RMSE - root mean squared error

As can be seen from the above overview, the DM performed best when projecting the data to quadratic space, while the GPC did best on the original data as-is. This was taken into account for the plots below for the DM and GP respectively, which used an instance of the more favourably performing processed data. Note that the exact probabilities provided by the GP are shown in the following plots, in contrast to the argmax taken for error-calculation purposes.

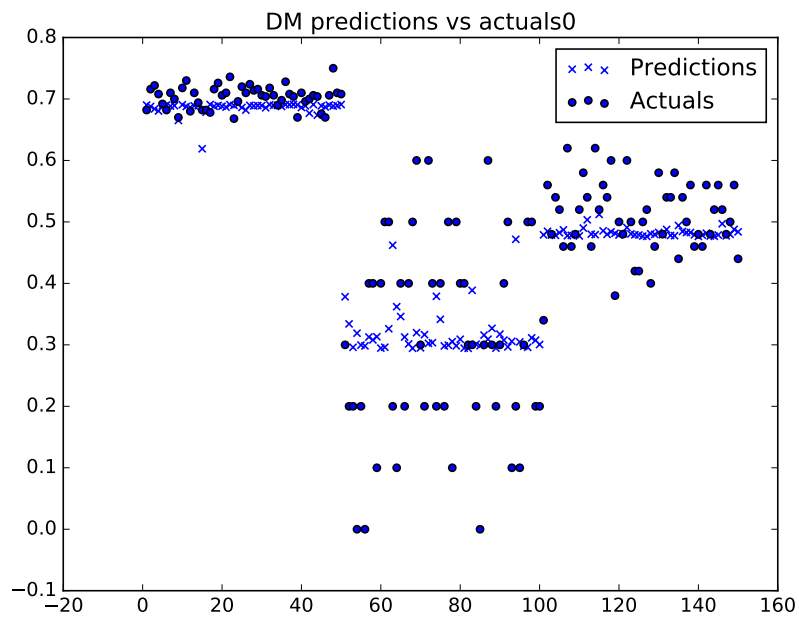


FIGURE 3.14: DM Label distribution of label 0

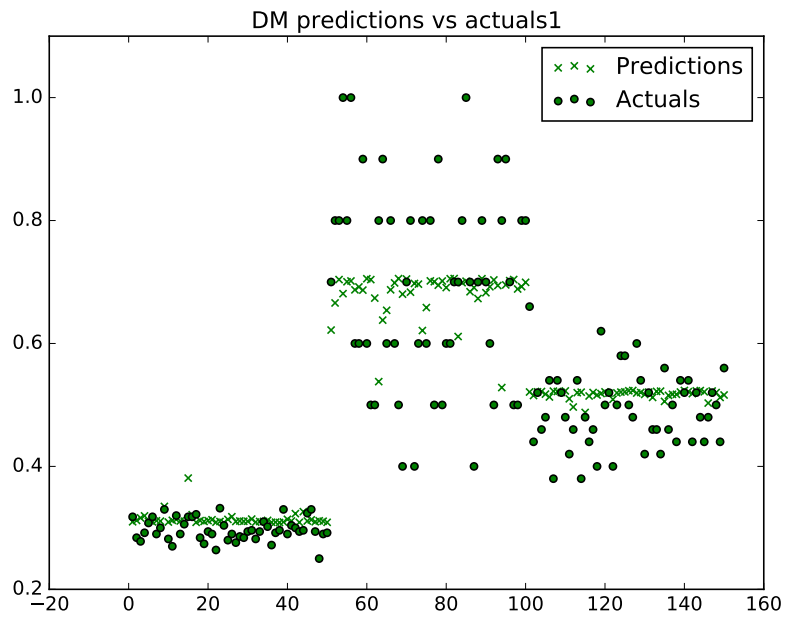


FIGURE 3.15: DM Label distribution of label 1

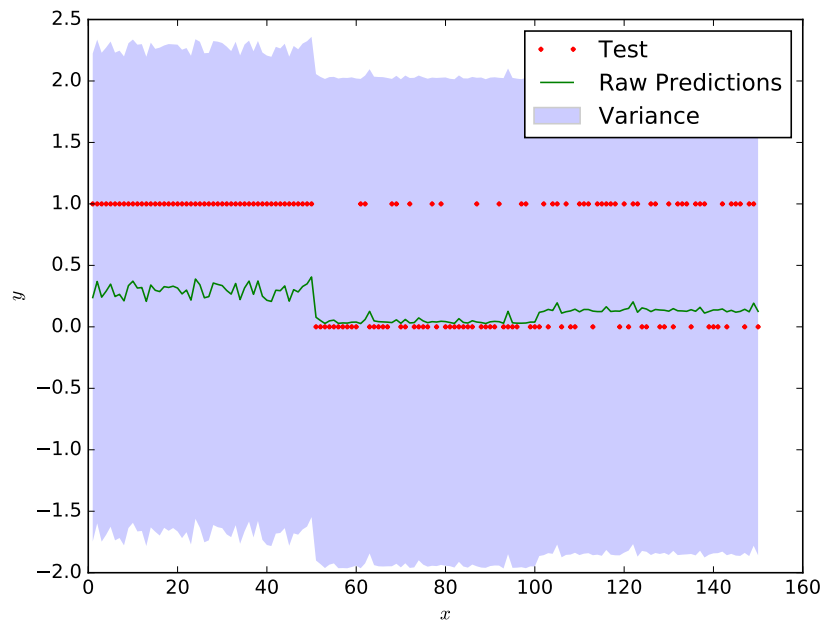


FIGURE 3.16: OvR GP performance with variance for label 0

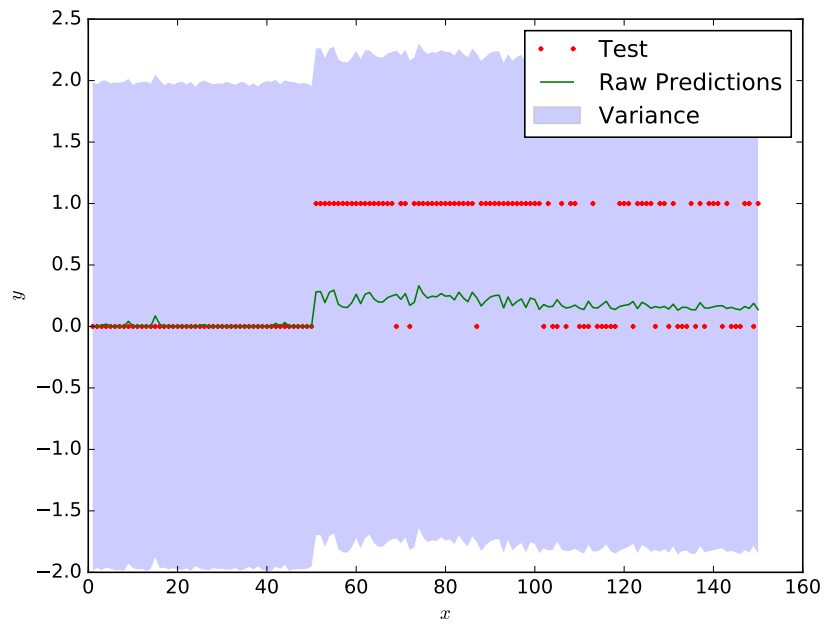


FIGURE 3.17: OvR GP performance with variance for label 1

As we can see, the DM performed notably better than the GP, though this is admittedly a rather simple example that assumes we have a sufficient amount of data from the *three* possible habitat clusters - A by itself, B by itself, and a homogeneous mix of A and B.

From this basic example, it is apparent that in the area where there is an even mix of labels A, B, the Gaussian Process' predictions are both noisy and very uncertain about their predictions, where human intervention would be required to observe the fact that it is in fact a consistent mix of both. In contrast, the dirichlet multinomial regressor is more confident in the fact that that area does in fact have a mix of labels.

CHAPTER 4

Results

The results from the Experiments detailed in Chapter 4 are listed below. The range of possible class values in some cases have been stretched beyond the existing class labels so that values align between different outputs to allow for easy, direct visual comparison. Note that the results to the above experiments will include those of both non-downsampled and downsampled results, as well as the full set of 24 labels as well as simplified ones.

Due to the low occurrence of some labels in the original dataset though, they have ended up being omitted in predictions - these are excluded from the colour schemes of the benthic maps generated, so that those that do occur can be given more distinct colours from one another as to better differentiate between the habitats of a map, as well as allow a consistent comparison of across different maps.

In this section, the performance of common machine learning algorithms, namely kNN, Logistic Regression, Random Forest, and SVM are explored first, to provide a comparison to the later, more complex algorithms.

4.1 Deterministic Methods

results here currently outperform the argmax of GPs and DMs...hmm

Algorithm	10F-CV F1	10F-CV Accuracy	Parameters	Data
KNN	0.7180212553324016	0.8752880347998457	n = 5	simple labels
Logistic Regression	0.21221134330881872	0.7374253253308372		simple labels
Random Forest	0.8103659316744121	0.9117076519281244		simple labels
SVC	0.21221261743619468	0.7374255088743278	OvA	simple labels
DM	0.287405310254214	0.757925654489819		simple labels
KNN	0.4711662374785644	0.6510716003156948	n = 5	full labels
Logistic Regression	0.06457891531653175	0.25912130389295746		full labels
Random Forest	0.6075895500970125	0.7263355175008718		full labels
SVC	0.012413092166172946	0.16549499843988033	OvA	full labels
DM	0.13802716811804644	0.37856057852908254		full labels

4.2 Gaussian Process Classification

transfer all the results from markdown

show more stratified results (not just even) to show that even did better

500	Even	GP	10F-CV	10	0.86534	Deterministic	
500	Stratified	GP	10F-CV	10	0.80136		
1000	Even	GP	All	1	0.87626		0.56208
1000	Even	PoEGP	All	5	0.80973		0.47481
1000	Even	PoEGP	All	200	0.80186		0.47595
1000	Even	GPoEGP	All	5	0.80864		0.51018
1000	Even	GPoEGP	All	200	0.80105		0.47748
1000	Even	BCM	All	5	0.80682		0.48167
1000	Even	BCM	All	200	0.80421	RBF, EP (default)	0.48227
1000	Even	GPy	All	1	0.87638		0.57013

highlight areas with low certainty, etc.

4.3 Dirichlet Multinomial Regression

transfer all the results from markdown

highlight areas with biodiversity, etc.

CHAPTER 5

Evaluation and Discussion

old note - these will likely be in results/described in experiments Looking at distributions of the GPs, dirichlet multinomial draws and how they perform beyond just taking argmax, etc. should go here

5.0.1 Limitations

- training data doesn't explore any particular area exhaustively - hard to verify how accurate any model is even if cross validation scores are high

CHAPTER 6

Conclusion

The conclusion goes here.

6.1 Future Work

- perform similar experiments on incrementally changing data every few years - observe biodiversity/habitat changes
- replace the simple activation function in the dirichlet multinomial with a more complex model like a GP
- previous work has been done for finding least certain areas of a GP to decide where to send AUV's to maximise resulting confidence in habitat labels - use entropy to be able to do the same with dirichlet multinomials, whilst overcoming the problem of areas with consistent heterogeneous labels that otherwise confuse GPs

Bibliography

- Nasir Ahsan, Stefan B. Williams, and Oscar Pizarro. 2011. Robust broad-scale benthic habitat mapping when training data is scarce.
- Christina L. Belanger, David Jablonski, Kaustuv Roy, Sarah K. Berke, Andrew Z. Krug, and James W. Valentine. 2012. Global environment predictors of benthic marine biogeographic structure. *Proceedings of the National Academy of Sciences*, 109.
- Asher Bender, Stefan B., Williams, and Oscar Pizarro. 2012. Classification with probabilistic targets.
- Peter J. Bickel and Elizaveta Levina. 2008. Regularized estimation of large covariance matrices. *The Annals of Statistics*, pages 199–227.
- Craig Brown, Stephen J Smith, and Peter Lawton. 2011. Benthic habitat mapping: A review of progress towards improved understanding of the spatial ecology of the seafloor using acoustic techniques. *Estuarine, Coastal and Shelf Science*, 92.
- J. Calvert, J.A. Strong, C. McGonigle, and R. Quinn. 2015. An evaluation of supervised and unsupervised classification techniques for marine benthic habitat mapping using multibeam echosounder data. *ICES Journal of Marine Science: Journal du Conseil*, 72:1498–1513.
- Rich Caruana and Alexandru Niculescu-Mizil. 2006. An empirical comparison of supervised learning algorithms. In *Proceedings of the 23rd International conference on Machine learning*. Association for Computer Machinery.
- Chang, Chih-Chung, and Chih-Jen Lin. 2011. Libsvm: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27.
- Robert J. Diaz, Martin Solan, and Raymond M. Valente. 2004. A review of approaches for classifying benthic habitats and evaluating habitat quality. *Journal of Environmental Management*, 73:161–181.
- Rozaimi Che Hasan, Daniel Ierodiaconou, Laurie Laurenson, and Alexandre Schimel. 2014. Integrating multibeam backscatter angular response, mosaic and bathymetry data for benthic habitat mapping. *PLoS ONE*, 9.
- Vladimir E. Kostylev, Brian J. Todd, Gordon B. J. Fader, R. C. Courtney, Gordon D. M. Cameron, and Richard A. Pickrill. 2001. Benthic habitat mapping on the scotian shelf based on multibeam bathymetry, surficial geology and sea floor photographs. *Marine Ecology Progress Series*, 219:121–137.
- Vanessa Lucieera, Nicole A. Hilla, Neville S. Barretta, and Scott Nichol. 2013. Do marine substrates “look” and “sound” the same? supervised classification of multibeam acoustic data using autonomous underwater vehicle images. *Estuarine, Coastal and Shelf Science*, 117:94–106.

- Arman Melkumyan and Fabio Ramos. 2009. A sparse covariance function for exact gaussian process inference in large datasets. *International Joint Conference on Artificial Intelligence*, 9.
- Aaron Micallef, Timothy P. Le Bas, Veerle A.I. Huvenne, Philippe Blondel, Veit Huhnerbach, and Alan Deidun. 2012. A multi-method approach for benthic habitat mapping of shallow costal areas with high-resolution multibeam data. *Continental Shelf Research*, pages 14–26.
- Kevin P. Murphy. 2012. *Machine Learning: A Probabilistic Perspective*. The MIT Press.
- Andrew Y. Ng and Michael I. Jordan. 2002. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. *Advances in neural information processing systems*, 14:841.
- Food & Agriculture Organisation of the United Nations. 2004. *The State of World Fisheries and Aquaculture*.
- Joaquin Quinonero-Candela and Carl Edward Rasmussen. 2005. A unifying view of sparse approximate gaussian process regression. *The Journal of Machine Learning Research*, 6:1939–1959.
- Carl Edward Rasmussen and Christopher K. I. Williams. 2006. *Gaussian Processes for Machine Learning*. The MIT Press.
- A. Rattray, D. Ierodiaconou, L. J. B. Laurenson J. Monk, and P. Kennedy. 2014. Quantification of spatial and thematic uncertainty in the application of underwater video for benthic habitat mapping. *Marine Geodesy*, 37:315–336.
- Jan Seiler, Ariell Friedman, Daniel Steinberg, Neville Barrett, Alan Williams, and Neil J. Holbrook. 2012. Image-based continental shelf habitat mapping using novel automated data extraction techniques. *Continental Shelf Research*, 45:87–97.
- Paul V R Snelgrove. 1994. Animal-sediment relationships revisited: Cause versus effect. *Oceanography and marine biology*.
- D. Steinberg, A. Friedman, O. Pizarro, Williams, and S.B. 2011. A bayesian nonparametric approach to clustering data from underwater robotic surveys. International Symposium on Robotics Research.