



A review of machine learning applications to coastal sediment transport and morphodynamics

Evan B. Goldstein^{a,*}, Giovanni Coco^b, Nathaniel G. Plant^c

^a Department of Geography, Environment, and Sustainability, University of North Carolina at Greensboro, Graham Building, 1009 Spring Garden St., Greensboro, NC 27412, USA

^b School of Environment, University of Auckland, New Zealand

^c U.S. Geological Survey St. Petersburg Coastal and Marine Science Center, St. Petersburg, FL, USA

ABSTRACT

A range of computer science methods termed machine learning (ML) enables the extraction of insight and quantitative relationships from multidimensional datasets. Here, we review the use of ML on supervised regression tasks in studies of coastal morphodynamics and sediment transport. We examine aspects of ‘what’ and ‘why’, such as ‘what’ science problems ML tools have been used to address, ‘what’ was learned when using ML, and ‘why’ authors used ML methods. We find a variety of research questions have been addressed, ranging from small-scale predictions of sediment transport to larger-scale sand bar morphodynamics and coastal overwash on a developed island. We find various reasons justify the use of ML, including maximize predictability, emulation of model components, the need for smooth and continuous nonlinear regression, and explicit inclusion of uncertainty. The expanding use of ML has allowed for an expanding set of questions to be addressed. After reviewing the studies we outline a set of best practices for coastal researchers using machine learning methods. Finally we suggest possible areas for future research, including the use of novel machine learning techniques and exploring open data that is becoming increasingly available.

1. Introduction

The amount of available data on coastal systems has increased in recent years, ranging from topographic and bathymetric data (e.g., Turner et al., 2016), to compilations of sediment transport and physical forcing (e.g., Bolaños and Souza, 2010; Garel and Ferreira, 2015; Nelson et al., 2013; van der Werf et al., 2009). Large spatial and temporal extents, high resolution, and rapid turnaround from acquisition to availability means that the data being produced enables expanded applications to coastal morphodynamic research. In particular, since observational data has always been the foundation for developing empirical relationships or testing quantitative models, the recent volume of data available, the intrinsic high dimensionality and nonlinearity of underlying processes (e.g., Rubin, 1992; Jaffe and Rubin, 1996; Werner, 1999; Murray et al., 2009; Murray et al., 2014a), and increased computing power, have all led to renewed interest in empirical research (e.g., Turner et al., 2016; Luijendijk et al., 2018).

A key example is attempting to extract insights, predictions, or quantitative relationships directly from multidimensional datasets using automated tools. This data-driven route for science has been demonstrated to be a promising research direction (e.g., Anderson, 2008; Hey et al., 2009), and tools from a range of disciplines have been influential in defining and tackling data-driven science (e.g., Stalzer and Mentzel, 2016). We differentiate classic empirical work and this

wave of data-driven work as being divided by the computational methodology, as well as (potentially) the quantity, nonlinearity expressed in the data, and high dimensionality. Our focus is on empirical work using machine learning, the set of computer algorithms, methods and tools that implement a given task and use data to optimize performance (e.g., reduction of error). In this manuscript we provide many examples of successes in the use of machine learning for coastal research, but first we discuss the rationale for this data-driven approach.

Data-driven research is inductive. As with other empirical work, data-driven research relies on data to develop insight, predictions, or relationships. Empirical work does not and cannot exist in a vacuum — theory and logic are critical parts of data analysis (e.g., Coveney et al., 2016; Crutchfield, 2014) and mathematical proofs show the lack of generalizability of inductive statements (e.g., Popper and Miller, 1983). However, inductive statements are part of the scientific workflow, and are unavoidable at certain junctures. Even Newton expressed the utility of induction in Rule 4 in the 3rd edition of the Principia (Cohen et al., 2016):

“In experimental philosophy, propositions gathered from phenomena by induction should be considered either exactly or very nearly true notwithstanding any contrary hypotheses, until yet other phenomena make such propositions either more exact or liable to exceptions.”

Coastal morphodynamics specifically and geomorphology in general

* Corresponding author.

E-mail address: ebgoldst@uncg.edu (E.B. Goldstein).

have long histories of induction, and developing empirical rules that are useful. Even when basic laws of physics can be used, empirical expressions or heuristic rules are always required to close an equation set. For example, sediment transport rules, turbulent closure schemes, friction coefficients, and wave breaking all rely on ad hoc rules, assumptions or empirical relations. If we still need inductive rules, how should we build them? With the increased quantity of data, and the improvements in computing power, coastal researchers have access to computer science tools from the subdiscipline of machine learning (ML) to develop inductive statements and optimized predictions directly from data sets. Machine learning differs from statistical learning since in ML there is no assumption or hypothesis about the structure of the relationship in the data, and instead there is an automated searching for rules and relationships. Also, in ML no restrictive assumption about the data is made (e.g., residuals do not need to follow a specific distribution). One of the arguments invoked to explain the better performance of ML predictors over statistical learning is in fact the lack of constraints related to data and assumptions.

Much of the machine learning work we discuss here is focused on identifying and exploiting correlations and patterns in data. Assigning causation can be less clear in some coastal morphodynamic systems because of multiple scales (ripples, megadunes, bars, shoreline), and feedbacks between scales (bars impact shoreline, and vice versa) that interact in both space and time (e.g. Murray et al., 2014a,b; Sherman, 1995; Wright et al. (1985); Werner, 1999; Winant et al., 1975). However, correlation is valuable for prediction because of the concept of analogy (Lorenz, 1969a,b) — knowing how a coastal system evolved when it was in the same configuration but at a previous time can lead to predictions about how the system might evolve in the future (i.e., seasonal dynamics of sediment transport; Aubrey, 1979; Plant et al., 1999, 2006; Splinter et al., 2011; Yates et al., 2009). A data-driven approach can help to elucidate this behavior by examining previously collected data, and developing a model focused on how the system will evolve based on past instances.

Data-driven work may only be strictly applicable within the range of the data used to develop the predictor — unless the prediction scheme can be argued to be more generally valid. This is a limit of all inductive, empirical techniques, though it is rarely mentioned in more traditional empirical studies (i.e., any study that uses linear regression to predict beyond the bounds of the data). Furthermore, this caveat is likely applicable to all modeling studies because processes or feedbacks can exist that are not included within the model. A morphodynamicist might argue that, if a given model is built from conservation laws, the model should be able to predict outside of the range of conditions where the model has been tested. This argument also holds for data-driven work — the data used to construct the model adheres to conservation laws, therefore predictions (built directly from data) might also adhere to these physical constraints outside of the range of data used to build the model. This is not an excuse to avoid re-appraisal of data-driven work when more observations are made.

While there are a number of aspects of coastal science that can and do benefit from machine learning, we focus here on predictions of coastal sediment transport, coastal morphology, and coastal morphodynamics. We focused on supervised learning, specifically regression tasks (predicting continuous dependent variables) and do not address classification tasks (predicting discrete dependent variables). Supervised learning involves input and output data that are linked (such as wave forcing and a given morphological configuration) with the goal of developing a function to relate the input to a corresponding output and emulating physical processes relationships that are either poorly understood or complex and difficult to capture with deterministic models (note that ML is also used for unsupervised tasks). Excluded here are prediction of forcing and fluid phenomena when no reference to sediment transport is given or studies focused on engineering and structures.

Previous ML work written for an Earth and Environmental Science

audience has focused primarily on introducing ML algorithms, in the form of both books (Hsieh, 2009) and papers (Chau, 2006; Hsieh, 2004; Valentine and Kalnins, 2016). This previous work serves a key role in connecting coastal scientists to ML tools (e.g., Jones and Maccarone, 2013). Our work intends to move the purview of these previous reviews, which focused on introducing ML algorithms, explaining ML algorithms, and providing a few select examples to review. In this document, we do provide a brief introduction to the most common machine learning techniques that have been used in coastal sediment transport and the steps needed to start a ML project (Section 2). However, our focus is on comprehensively reviewing previous machine learning work on coastal morphodynamics such that these works can be recognized, compared, and used to build future ML efforts. We review and discuss more than 60 papers under three separate headings: studies where ML is used to predict sediment transport (Section 3); studies where ML is used to make a morphodynamic model (Section 4); and studies where ML is embedded in or linked to a morphodynamic model (Section 5). Furthermore, we address why ML tools have been used in particular studies and what was learned by using ML methods. Coastal morphodynamicists interested in ML can use these sections (3, 4, and 5) to assess what research has been done (and what has not been done) at the time of this writing. This comprehensive review of the coastal ML literature also permits us (in Section 6) to discuss overarching topics, offer a set of best practices for open, reproducible machine learning research and highlighting some future directions in coastal ML research primarily focused on synthesis and intercomparison.

2. Machine learning methods used in coastal sediment transport and morphodynamics

Before we review the uses of ML in coastal science, we briefly introduce the ML methods such that sections (Sections 2.1–2.4) provide basic information on each ML method that has been used in the reviewed coastal works. There are many additional ML techniques that are not discussed because we could not find examples of their use in coastal morphodynamics and coastal sediment transport research. Within each of these sections we provide relevant papers for readers who wish for more details regarding each method

2.1. Artificial neural network (ANN)

Artificial Neural Networks (ANN) are commonly used algorithms in machine learning because of their versatility. Many different fields of science have used ANNs for tasks such as function fitting to classification. Applications in coastal sediment transport and morphodynamics include multiple aspects of suspended sediment transport, sandbar morphodynamics, and various studies of shoreline position — all mentioned in the following sections. The most typical form of an ANN used in coastal research is a multilayer perceptron, which is represented by a series of layers: an input layer, one or more hidden layers and an output layer. Each layer consists of a number of nodes (artificial neurons). The input data is fed to the network via a node on the input layer (usually each node represents an input variable). The hidden layer(s) contain a somewhat arbitrary number of nodes chosen based on a mix between experience, empirical formulas and systematic analysis. At the end of the network, depending on the number of variables to be predicted, the output layer could consist of one or more nodes. As information passes through the network (from input to hidden layer to output), it is modified at each node by a transfer function, which introduces nonlinearity into the ANN. An idealized feed-forward ANN characterized by n input nodes (the predictors or independent variables), m hidden nodes and one output node (the prediction or dependent variable). Nodes are mathematically connected and transfer information from the input variables to a node of the hidden layer:

$$h_j = f \left(a_j + \sum_{i=1}^n w_i x_i \right) \quad (1)$$

where x_i is the i th of n input variables, h_j the response of the j th neuron in the hidden layer, f is the transfer function (e.g., a sigmoid, an hyperbolic tangent, etc.), w_i is the connection weight between x_i and h_j , a_j is the bias for the j th hidden neuron. A further combination of the hidden nodes, which is achieved by means of a new activation function (not necessarily the same as the one used to link the input variables and the hidden layer) and new connection weights and biases, connects the hidden layer to the output layer.

The biases and connection weights of the ANN are established through an optimization algorithm that is applied to a dataset consisting of observed input and output variables. Various algorithms can be used to perform this critical step, though it remains difficult to tell *a priori* which optimization function will provide better results. Many of these algorithms are based on the backpropagation of errors — the error at the output (prediction) layer is sent back through the network to adjust and update the weights and biases.

ANNs are often portrayed as an example of a black-box predictor where the (usually) large number of weights and biases obscures the role of individual variables. The architecture of small-size ANNs can in fact be analyzed and various techniques have been developed to this aim (e.g., Olden et al., 2004). LeCun et al. (2015) provide many helpful references and a relatively recent review on ANN that focuses on current research themes (i.e., Deep Learning; ANNs with many hidden layers).

2.2. Genetic algorithms (GA) and genetic programming (GP)

Genetic algorithms (GA; Holland, 1975) and genetic programming (GP; Koza, 1992) are related ML techniques that operate on rules based on natural selection. In the section below we review the basics of genetic algorithms. For example, consider an equation with five free parameters, where a given combination of specific values for each parameter can be compressed into a single vector of length 5. The vector of parameter values is also related to the solution of the equation using these 5 specific values. Each solution is also related to an associated error (the value of the equation using the 5 parameters vs. some measured value). Now consider a population of such vectors (not just one), each with their own unique combination of values for each of the 5 parameters. The genetic algorithm routine works by operating on these vectors using evolutionary rules. Given an initial set of vectors (a population), there is an error associated for each vector. The vectors with the smallest error are retained; the vectors with the most error are discarded. New vectors are developed by mutation (changing values in a given vector) and by reproducing — recombining two vectors to make a new vector. By using these evolutionary rules, the routine will search over the solution space and tend to converge on solutions that are globally optimal. Parameters in the evolutionary rules and techniques in applying those rules are tunable (e.g., number of predictors in the population, mutation percent, crossover rules for combining parts of vectors with each other, number of generations, number of discarded or kept predictors for each generation, etc.). Genetic algorithms can be used in tandem with artificial neural networks — for example to find the appropriate weights and biases, as well as network architecture (e.g., Yao, 1999). Further helpful entries into the GA literature can be found in D'Ambrosio et al. (2013), Mitchell (1995) and Mitchell (1998).

Building on the population approach of genetic algorithms, genetic programming (GP; Koza, 1992) takes the idea a step further. The population is not vectors of parameter values to be input into a fixed equation but instead a population of equations with mutable form and length. Given a set of mathematical operators (+, -, *, /), and a set of input variables (e.g., forcing conditions) a GP routine works to find equations using these building blocks (input variables, constants, and

mathematical functions) — this is a symbolic regression problem. One issue with GP is the development of large, complex functions that have smaller error compared to small, less complex functions that have larger error but might be more physically interpretable. Therefore routines may offer more than one solution, and instead offer many solutions to the problem which fall along the Pareto front — a line in error-complexity space that defines how prediction error decreases with the solution complexity (a measure of the size of the predictor that incorporates the mathematical operators, variables, and constants). The act of choosing a predictor from this front introduces subjectivity in the routine, though GP algorithms have shown the ability to find physically meaningful results from data streams. Aside from the work of Koza (1992) introducing the technique, the book by Poli et al. (2008), and work by Babovic and Keijzer (2000), Olden et al. (2008), Schmidt and Lipson (2009), and O'Neill et al. (2010) have proven helpful to us.

2.3. Bayesian networks (BN)

Bayesian networks (BN) implement a form of probabilistic prediction that explicitly resolves the conditional probabilities that link variables to one another, albeit in a discretized fashion. Statistical operations include marginalization over a subset of a larger distribution, for instance, when the data are used to provide constraints (Charniak, 1991). And, as the name suggests, Bayesian estimation (Cooper and Herskovits, 1992; Malakoff, 1999) can be implemented to solve problems that typically require data assimilation (Wikle and Berliner, 2007). For example, to estimate coastal erosion that is assumed to be influenced by dune morphology, geology, and sea-level rise (Plant et al., 2016) the probabilistic relationship can be expressed as:

$$P(E_i) = \sum_{G,D,SLR} P(E_i | D, G, SLR) P(D | G, SLR) P(G | SLR) P(D) P(SLR), \quad (2)$$

where left side of the equation is describing the probability that a certain amount of erosion, E_i , is experienced. The right side of the equation is the product of the conditional probability of that amount of erosion occurring, given the morphologic state of coastal dunes (D) and geologic setting (G) and a sea-level rise rate (SLR). The probability is integrated over all the states, which may be constrained by data. This is the marginalization operation that assumes conditional probabilities have been estimated via a learning process. Some of the terms on the right side of Eq. (2) defining the erosion probability may themselves have dependencies that can be solved using Bayes rule:

$$P(D_i) = \sum_E P(E | D_i, G, SLR) P(D_i) / P(E) \quad (3)$$

where the first term on the right side of the Eq. (2) is inverted. Bayes rule and marginalization can take place simultaneously in a Bayesian network, implying that there is no real distinction between a forward implementation that emulates a deterministic model (e.g., a partial differential equation) and an inverse model.

The approach models probabilities directly, as opposed to modeling the process-variables as is done in the other ML examples. This is useful if knowing the uncertainties is a primary modeling requirement. A disadvantage is that the model must learn the conditional probabilities that describe the correlations between variables and this comes with a cost of increasing free parameters that grows as the number of states raised to the number of variables. Furthermore the uncertainty present in the resulting model only reflects the uncertainty that is found within the data. We have found general papers by Aguilera et al. (2011), Chen and Pollino (2012), Uusitalo (2007), Beuzen et al. (2018b), and Beuzen and Simmons (2019) to be useful in learning techniques and applications of BN. Other ML techniques we discuss in this review can be used to make probabilistic predictions (e.g., by developing many ANN predictors for the same problem while using different partitions of the training data), and there are many other probabilistic ML techniques

that exist but remain unused by coastal scientists.

2.4. Regression trees (RT)

Regression trees (RT) separate prediction tasks into a series of binary splits, leading to a branching, tree-like structure (e.g., De'ath and Fabricius, 2000; Hastie et al., 2009). Tree based approaches are often used for classification task because easily allow users to assess the relative influence of the input variables. Binary splits and discrete output given specific input conditions are not ideal when dealing with continuous dependent variables that vary over a large range, but they have been used in relevant coastal research. Trees are visually appealing and reading through a RT model can be straightforward, especially when the tree is short. Regression trees also have a tendency to be unstable to changes in training data, developing overfit models and performing poorly at generalization if users do not compensate (e.g., via 'pruning')

An example of a regression tree-based algorithm is recursively splitting the dataset into groups. The details of each split are determined via a given metric, such as minimizing the sum of squares of each group. A variety of rules exist for both growing trees (i.e., how many recursive splits) and pruning trees (removing splits). Additionally, other algorithms can be attached to tree based methods to improve accuracy, specifically Boosting. Boosting routines merge many small regression tree models that are built sequentially, with misfit data sequentially given more weights so trees progressively focus on poorly predicted data (Elith et al., 2008). We have found the works of De'ath and Fabricius (2000), De'ath (2007), Olden et al. (2008), and Hastie et al. (2009), useful for learning about these tree based approaches.

3. Applications to coastal sediment transport

We have now reviewed the most used ML techniques in coastal morphodynamics and sediment transport studies. The availability of coastal sediment transport data, and the lack of a single perfect predictor (for a given sediment transport relation) has lead to the hope that ML will provide a more viable, optimal sediment transport equation — a motivation for many of the works that we review in this section. Authors frequently want to develop a predictor that is either more generally valid (better prediction with a large set of data) or more specifically valid (better prediction with a small set of data specifically collected for a given setting/condition). Authors of the studies reviewed below all test their ML prediction scheme against established predictors from the literature (i.e., previous empirical or theoretical sediment transport prediction schemes). The newly developed ML techniques often performs better than the traditional scheme using the error metric selected by the authors, a phenomena we discuss in Section 6.1.

3.1. Suspended sediment concentration

Predictions of suspended sediment concentration are a fundamental test of theoretical and statistical understanding of sediment mobility and transport that control morphologic evolution on a range of spatial and temporal scales. Time-varying sediment concentrations have been predicted using several ML methods. Jaffe and Rubin (1996) used nonlinear forecasting techniques to predict suspended sediment concentration based on instantaneous water velocity (with and without higher order velocity terms) and various water velocity history terms (e.g., velocity at the previous time step, etc.). A notable aspect of this study is the investigation of the appropriate time lag in water velocity to maximize the correlation with sediment transport (Fig. 1) an insight that may have transferability to other studies on time lags in coastal systems. As an extension to this work, Yoon et al. (2013) used an ANN to predict time-dependent suspended sediment concentration as a function of various hydrodynamic parameters both inside and outside the surf zone. With such a large dataset and many measured variables, Yoon et al. (2013) was able to use the ANN to identify the

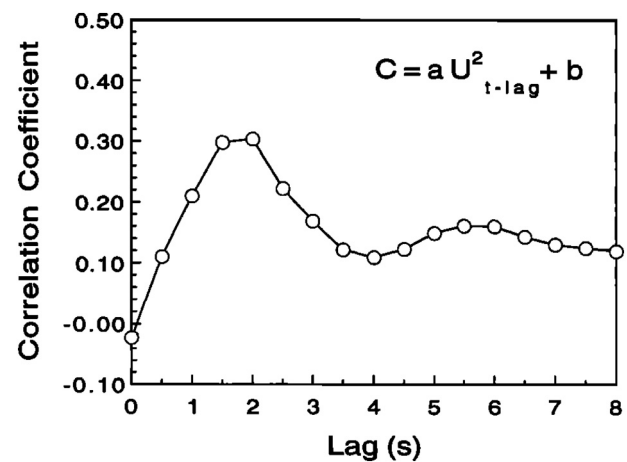


Fig. 1. A plot from Jaffe and Rubin (1996), who used nonlinear forecasting to predict suspended sediment concentration under waves. The plot above exhibits the changing correlation between suspended sediment and wave forcing with changing lag time on the wave forcing term.

hydrodynamic parameters (and combinations of parameters) that are most predictive in different regions of the laboratory surfzone. Using a GP routine, Kizhisseri et al. (2005) used both synthetic and field data to produce expressions for suspended sediment concentration based on instantaneous fluid velocity (and higher powers of velocity). In a rare example of a reported unsuccessful ML application in coastal morphodynamics, the prediction of suspended sediment concentration using field data lead to poor performance (i.e., a large absolute error for the prediction; Kizhisseri et al., 2005). Oehler et al. (2011) used both ANN and Boosted Regression Trees (BRT) to develop predictors for near bed suspended sediment reference concentration based on water depth, median grain size, mean wave period at the bed, wave orbital amplitude at the bed, and significant wave orbital speed at the bed. The BRT model was superior to ANN (Fig. 2), which we highlight because many studies do not compare ML derived predictors developed from multiple ML routines. Oehler et al. (2011) provides a clear example that this work should be done, and will inform future researchers interested in which ML method to use to predict suspended sediment reference concentration. Goldstein and Coco (2014) used the same dataset and developed a GP routine to construct a predictor for reference concentration. This predictor was specifically derived for use in a numerical model of inner shelf bedforms (discussed further in Section 5), and is an example of a predictor developed to work in a specific (multiple grain size) setting.

3.2. Suspended sediment flux

Scaling up from instantaneous concentration to alongshore-directed suspended sediment flux has been the focus of several studies. Using an ANN, van Maanen et al. (2010) predicted the depth integrated alongshore sediment transport using water depth, wave height, wave period and alongshore current velocity. Analyzing the parameterized ANN also allowed van Maanen et al. (2010) to understand which parameters held the most explanatory power (alongshore current of velocity), and to understand when the predictor provided unphysical answers. Notably, unphysical predictions were found when the ANN was given input parameters outside the range of the training data, highlighting the importance of training models with extreme conditions. Predictors for the net alongshore sediment transport rate based on wave height, wave period, breaking wave angle, beach slope, and grain size have been developed using ANN (Kabiri-Samani et al., 2011) and regression trees (Mafi et al., 2013). The ability for both ML methods to produce successful predictors highlights the need for more comparative work between ML methods.

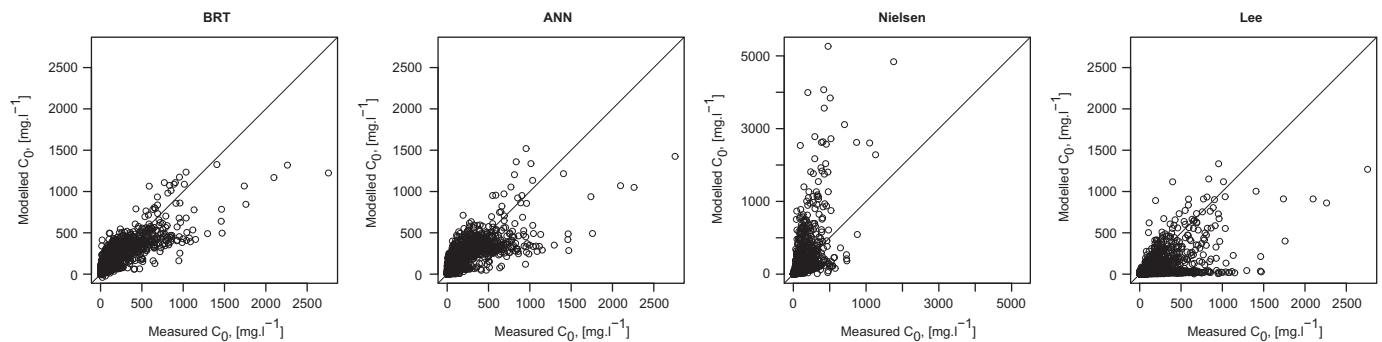


Fig. 2. A figure from Oehler et al. (2011) exhibiting the performance of a BRT and ANN model for suspended sediment reference concentration compared to two more traditional prediction schemes.

3.3. Sediment properties

Finally there have been ML studies of sediment properties (e.g., mean grain size, skewness, kurtosis, fall velocity, etc.). Nylén et al. (2015) trained a decision tree to determine several aspects of beach and dune sediment in Finland as a function of environmental variables (e.g., elevation, slope, curvature, local fetch, geography and climate conditions). Sediment parent material (parameterized via location) was found to be an important control on grain size and sediment sorting, obscuring the role of local controls. Goldstein and Coco (2014) used a GP routine to develop a predictor for noncohesive sediment settling velocity that incorporates fluid kinematic viscosity, relative sediment density and sediment nominal diameter. The study focused on the role of training dataset size and selection method while developing a prediction scheme that performed better than two common equations.

4. ML morphological and morphodynamic models

A variety of coastal morphology and morphodynamic models have been built using ML. Many researchers use ML as an optimization tool — looking for better morphological prediction with newly collected or existing data.

4.1. Sandbars

Sandbar morphology (e.g., the cross-shore position and alongshore uniformity) has been a common focus of machine learning studies. Múnera et al. (2014) developed an ANN to determine the correlation between sandbar morphology and a given wave climate, culminating in examining the nonlinear dependencies of bar position on past wave conditions (i.e., time-lagged wave conditions). López et al. (2017) used an ANN to determine the cross-shore bar position given wave characteristics, sediment characteristics, and temporal data (month and day information). Compared to a common formula to predict bar characteristics, the optimized ANN had lower error. Komurcu et al. (2013) used an ANN to predict the geometric and shape characteristics of experimentally simulated bars based on the wave height, wave period, bed slope and grain size. Tests were performed varying the split of training data/testing data. The best fit model was trained with the largest amount of data, and had lower error compared to literature formula. A similar study on experimental bar data was performed by Demirci et al. (2015), using wave parameters, bed slope and sediment characteristics to predict bar volume using an ANN and multiple linear regression. The predictor derived from ANN outperformed the multiple linear regression.

Of particular note is the work of Pape et al. (2007, 2010), who used a recurrent artificial neural network to model the cross-shore position and temporal dynamics of sandbar crests. Recurrent neural networks are ANNs that feed output predictions back to the input layer of the ANN, making a forward in time morphodynamic model. Pape et al.

(2007) modeled sandbar position using relevant wave inputs and previous sandbar positions using a linear autoregressive model with exogenous inputs and a recurrent neural network (i.e., a nonlinear autoregressive model with exogenous inputs; NARX) trained using multiple techniques. All models exhibit decaying performance as the prediction horizon (the prediction lead time) increases, but nonlinear ANN models show slightly better results over long prediction timescales (Fig. 3). Assessment of prediction timescale is especially critical to understand if data-driven techniques can be successful techniques for forecasting future morphology and morphodynamics, and understanding how error compounds or decays through time in data-driven models. Additionally, comparative work between data-driven methods is particularly interesting to the user of ML, and can give insight into inherent predictability and nonlinearity in the study system. Pape et al. (2010) continued this work, using two neural networks to model sandbar behavior and compared results to a traditional cross-shore morphodynamic model. Both data-driven models showed performed better than the morphodynamic model, measured using the metric of error over increasing prediction horizon. Recurrent neural networks are more difficult to train than simple multilayer perceptron ANN, but significant advances in training methods and network architecture (see LeCun et al., 2015) have occurred since the work of Pape et al. (2010).

4.2. Shoreline position and shore profile

Various shoreline attributes have also been predicted using machine learning techniques. Using an ANN to predict beach profiles (from the dunes to MSL) with wind and wave data, Hashemi et al. (2010) also discusses the role of training data that spans a range of conditions to avoid error associated with out-of-sample prediction. Grimes et al. (2015) analyzed beachface and shoreline timeseries data, using a GP and nonlinear forecasting to predict the dynamics of intertidal

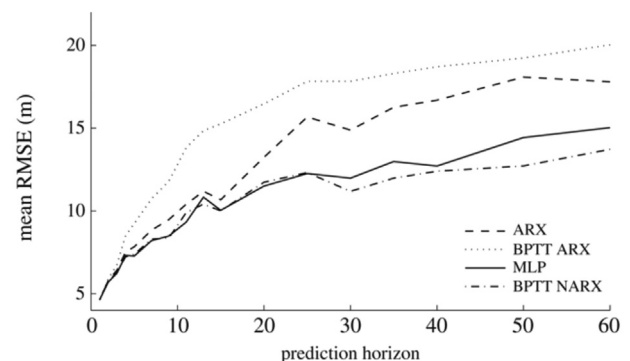


Fig. 3. The increase in sandbar position error for increasing prediction timescales using data-driven models (from Pape et al., 2007). Error is minimized when using models based on ANN (the ‘MLP’ and ‘BPTT NARX’ models).

beachface geometry and examine the role of internal dynamics vs. external controls (i.e., forcing). Both Tsekouras et al. (2015) and Rigos et al. (2016b) described new methods to formulate and train a novel ANN architecture to predict shoreline characteristics — Tsekouras et al. (2015) examined shoreline erosion as a function of storm characteristics and bathymetry, while Rigos et al. (2016b) investigated multiple shoreline positions and shoreline rotation given hydrodynamic inputs and offshore reef morphology. Iglesias et al. (2009a, 2010) used an ANN to predict the planform morphology of headland-bay-beach systems (including those with shore protection structures). Iglesias et al. (2009a) tested multiple ANN architectures (number of hidden layers and nodes) as well as different algorithms to train the ANN. The final ANN model outperformed previously developed shoreline models, and error from the ANN model was distributed across the shoreline as opposed to the previous models, which had concentrated zones of high error. Iglesias et al. (2009b) extended this work by incorporating tidal range into the ANN. After testing various ANN architectures and finding the best ANN predictor, Iglesias et al. (2009b) used the trained model to examine the interplay between tidal range and wave parameters in controlling headland bay geometry. Loureiro et al. (2013) use a BN to probabilistically determine the beach state classification (i.e., Wright and Short, 1984) given a range of hydrodynamic data and sedimentological data. The study found utility in the uncertainty of prediction, namely a range of possible beach states that are likely to exist at a given site based on environmental conditions. Bayesian networks have also been used to make probabilistic predictions of coastal morphology at large scales. Gutierrez et al. (2011) use a BN to develop shoreline change rate predictions for the US east coast based on hydrodynamic, simplified geomorphology (with bins based on vulnerability to sea level rise), and cross shore morphology of the coastline. Plant et al. (2016) modify the BN of Gutierrez et al. (2011) to include and probabilistically predict shoreline change as well as dune height for work in the Gulf of Mexico. Interestingly, the inclusion of dune height as an input variable increases precision but predictions are not more accurate. Yates and Le Cozannet (2012) also use the approach of Gutierrez et al. (2011) to probabilistically assess future European coastal evolution (either erosion, stable or accretion) using geomorphology (simplified coastal type; e.g., wetland, rocky cliff, beach, etc.), simplified geology (hard vs. soft sediments), mean tidal range, rate of sea level rise (RSLR), and mean significant wave height. Stable coastlines are predicted with greater accuracy compared to erosive or accreting coastlines, and the authors suggest that incorporating more local behaviors may resolve this issue. Bulteau et al. (2015) use a Bayesian network to predict shoreline change on La Réunion island from geomorphic setting (e.g., cliff, shingle beach, sand beach, etc.), presence of an estuary, anthropogenic structures, RSLR, and a function of wave energy. In addition to discussing why certain inputs are more predictive than others (specifically local geomorphology), this study also specifically examines areas of misprediction, and offers insight as to the unique situations when misprediction arises. Lentz et al. (2016) use a BN to relate land cover classifications, present elevations, and expected changes in RSLR to likelihood that coastal geomorphic settings would evolve to keep up with sea-level rise or inundate. A BN approach was used specifically in this study because of its probabilistic nature. Gutierrez et al. (2015) predicted decadal changes in barrier island geomorphology of Assateague Island, USA with a Bayesian network that uses a range of input parameters such as shoreline change rate, distance to an inlet and local morphological measurements (e.g., dune height, beach width) in a Bayesian network.

In addition to looking at long term shoreline changes, event scale work has also used BNs for prediction. Wilson et al. (2015) built on previous work by Lentz and Hapke (2011) to predict the beach volume changes resulting from storm events on Fire Island NY, USA with a BN. Predictions were improved in this network by anthropogenic impacts on the beach (nourishment in this location) and adjusting the hydrodynamic inputs to the model (runup elevation vs., impact hours). This

highlights the potential role of using several different inputs that may be viewed as quantifying the same process (wave-driven erosion), but may vary in correlation with the desired output (beach volume change). Beuzen et al. (2017) compared the use a BNs of different size (2 vs. 3 input nodes) to predict shoreline retreat as a consequence of storm events and preexisting beach characteristics (state, slope, width) at Collaroy-Narrabeen Beach in SE Australia. Building on this work, Beuzen et al. (2018a) examined the use of BNs as both predictive tools (high performance on testing data) and descriptive tools (high performance on training data) for storm-driven shoreline change. Beuzen et al. (2018a) notes that BNs built for descriptive purposes can be used to gain insight on underlying processes that produce the data, including causality.

Bayesian networks have also focused on emulating process-based models of storm erosion that are particularly computationally intensive—Poelhekke et al. (2016) use a BN as an emulator for the detailed process based model XBeach (Roelvink et al., 2009). By developing a set of forcing conditions for XBeach and running the model for each forcing condition, Poelhekke et al. (2016) train the BN to predict morphodynamic impacts (overwash depth, flow velocity, and erosion) on Praia de Faro, Portugal, a developed barrier island. The goal of the work is to develop a quick method to emulate the XBeach for use in an early warning system. Plomaritis et al. (2018) extended the work of Poelhekke et al. (2016) and a BN on disaster risk reduction (Jäger et al., 2018) to assess the impact of risk reduction measures on morphodynamic impacts on the Ria Formosa Barrier system of Portugal. Again the BN served as mechanism to emulate process based model runs.

4.3. Dune erosion

BN studies of coastal dune erosion from storm events have also found value in explicitly developing probabilistic prediction — Plant and Stockdon (2012) use a BN to predict dune crest elevation changes, dune crest position change, and shoreline position change as a function of dune base elevation, storm induced mean water level, and storm induced run-up (Fig. 4). Observations of dune erosion do not always match predictions perfectly, but do fall within the confidence intervals of the probabilistic method — this highlights the utility of probabilistic predictions toward enhancing prediction accuracy and certainty.

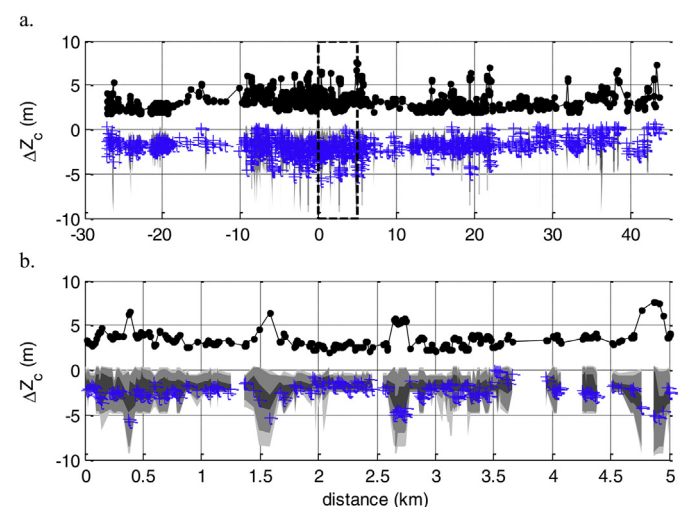


Fig. 4. Figure from Plant and Stockdon (2012), who used a Bayesian network to make predictions of foredune crest elevation change (ΔZ_c). Initial dune height is shown in black, observations are blue +, and predictions from the BN are shown as shaded area ranging from 50% (dark) to 95% (light) confidence interval. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

Palmsten et al. (2014) used the network from Plant and Stockdon (2012) as well as a simplified model structure to develop probabilistic predictions of dune position change along the Gold Coast in Queensland, Australia. Of note in this study is the attempt to use the trained model from Plant and Stockdon (2012), with no modifications or additional training, for a new site — prediction was not skillful with this model, however the ability for ML models to be generalized and extrapolated to new sites is an important test for any coastal ML models. den Heijer et al. (2012) performed a similar test using a Bayesian network that was designed to emulate an existing volumetric dune erosion model. The trained model was not able to successfully extrapolate beyond the range of the training data.

4.4. Cliffs and rocky coastlines

Much of the previous work has focused on low-sloped sandy coastlines, though there has been work on rocky coastlines. Dickson and Perry (2016) use several regression tree approaches to identify the controls on coastal cliff landsliding (e.g., distance to fault, bedding dip, aspect, etc.). Multiple methods converged on the same two controlling variables, a benefit when comparing multiple ML methods. Hapke and Plant (2010) use a BN to develop a relationship between short term cliff erosion rate of rocky coastlines of the southern California, US, and underlying geology, cliff height, cliff slope, and a metric based on hours the cliff is subject to wave attack, and long term erosion rate of the cliffs.

4.5. Wave ripples

Shifting from the coastline to smaller scale morphology, Yan and Zhang (2008) built an ANN to predict wave generated ripple size (length and height) based on sedimentological and hydrodynamic conditions. Data was both field and lab studies, and the ANN results were compared to four other empirical models. The ANN results provide more accurate predictions based on 3 statistical measures (scatter index, correlation coefficient, and mean geometric deviation) than the empirical models. Also studying wave-generated ripple geometry, Goldstein et al. (2013) used a GP routine to construct an equation for wave generated ripple height, wavelength, and steepness using sediment grain size and near bed orbital excursion. The new machine learning scheme produced more accurate predictions compared to traditional predictors. This predictor was ultimately used as a component within a larger numerical model (Goldstein et al., 2014).

4.6. Detection of bars and shoreline in images

Detection of morphological features from video images has also employed regression-based ML. Kingston et al. (2000) used an ANN to model the difference between sandbar position and video intensity maxima with additional inputs of wave height and tide level. Additionally, the model developed in Kingston et al. (2000) showed success against other methods (Plant et al., 2007). Related work has focused on detecting the shoreline in video observations with a variety of ANN architectures (Alvarez-Ellacuria et al., 2011; Rigos et al., 2016a).

5. Hybrid ML morphodynamic models

Machine learning methods can be linked with morphodynamic models to create what we refer to as hybrid models, after Krasnopolsky and Fox-Rabinovitz (2006), and Krasnopolsky (2013). There are several reasons for a hybrid models (Goldstein and Coco, 2015): 1) ML components can serve as emulations of complex routines or equations to speed up the computational process; 2) Data-driven parameterizations can serve as model components when parameterizations have ample data but no single optimal expression — perhaps there are multiple competing formulations; 3) More data might be anticipated in the near

term future, and the parameterizations might be volatile, subject to change as new data is collected; 4) Hybrid models offer a degree of specificity to a model. Adding a ML predictor is way of incorporating a bespoke prediction scheme, which can be useful for modeling a specific setting where data was collected.

Three coastal morphodynamic models have combined genetic programming routines to aid in various aspects of modeling. First, Goldstein et al. (2014) incorporated a GP derived suspended sediment reference concentration predictor and equilibrium wave orbital ripple morphology predictor (Goldstein et al., 2013) into a model of inner shelf sorted bedforms. The model previously had been built using theoretical and empirical parameterizations of these processes, but data from inner shelf sorted bedforms was used to develop new parameterizations and produce a refined model. The goal of the modeling work was to add more specificity to the process parameterizations in settings with mixed grain sizes. Second, Limber et al. (2014) and Limber and Murray (2014) used a GP derived expression as a component in nonlinear dynamical system model for rocky coastline evolution. The GP routine was used to develop an expression that emulated the output of a wave ray tracing model, thus summarizing the wave model results into a single smooth continuous equation amenable to further numerical work and phase plane analysis. Third, Goldstein and Moore (2016) developed a model of coastal dunes subject to storms that combines an empirical formulation of coastal foredune growth with a parameterization for dune erosion built using a GP. The GP routine was used to fit a smooth continuous equation to a set of data to facilitate numerical analysis.

Bayesian networks have been used as subcomponents for a variety of coastal models. Plant et al. (2014) used a Bayesian network to estimate overwash probability of a berm from hydrodynamic and wind conditions. This overwash probability was linked to a (non-BN) model of berm morphology. At larger space and time scales, both Passeri et al. (2016) and Bilskie et al. (2016) use the BN of Plant et al. (2016) — itself an extension of Gutierrez et al. (2014) — as a model component to predict century-scale shoreline change and dune height change as a function of SLR scenarios and geological constraints for the Gulf of Mexico. Bilskie et al. (2016) used the Plant et al. (2016) BN as a component in a larger model addressing Hurricane impacts under different SLR scenarios. Passeri et al. (2016) used the Plant et al. (2016) BN as a component of a model to simulate tidal hydrodynamics under SLR scenarios. Both Passeri et al. (2016) and Bilskie et al. (2016) mention that the Bayesian network was used because it is computationally efficient for the long time and large space scales that were modeled. Passeri et al. (2016) also discuss limitations to the BN component — the lack of historic data to train the BN limited its use in bays and estuaries, and large scale barrier island processes such as rollover of back barrier shoreline migration and nourishment also were not encoded in the BN (but included as rules in the larger model). Van Verseveld et al. (2015) used the process based model XBeach to simulate a storm event impacts on a developed barrier island. Hydrodynamic and morphodynamic output from the model was used as input for a Bayesian network, which predicted the damage to buildings. Van Verseveld et al. (2015) notes that BN require a significant amount of data. The use of multiple inputs (flooding, scour, wave height) and the probabilistic nature of the BN were also advantages in this research.

Many coastal morphodynamic models have several free parameters that must be tuned for a given field site or use case (e.g., Apotsos et al., 2008; Lin and Sheng, 2017; Murray et al., 2016; Pinsky et al., 2013; Plant and Holland, 2011; Stephens et al., 2011; Stockdon et al., 2014). When the number of free parameters is large and potentially inter-related, machine learning can be used to find optimal parameters. Knaapen and Hulscher (2002) developed a model for sand wave growth and saturated morphology, with best-fit model parameters are found using of a GA routine. Ruessink (2005) tuned nearshore model parameters using a genetic algorithm coupled to a local optimization routine (downhill simplex). Komurcu et al. (2008) used a GA to determine the

values for coefficients in two highly nonlinear functions that predict experimentally produced bar geometry based on the wave height, wave period, bed slope and grain size. Goldstein and Moore (2018) used a GA routine to tune a spatially explicit model of coastal foredune growth using structure-from-motion derived digital surface models.

Beyond using ML inside morphodynamic models (hybrid models) and to aid in model tuning, ML can also be used to analyze model results and gain insight into model output. For example, Lazarus et al. (2011) use nonlinear forecasting to quantify the nonlinearity of time-series output generated by a model of human interaction with the coastline. With the increasing complexity of numerical models, ML offers a way to analyze model results, leading to broader understanding applicability of numerical models and development of new hybrid models.

6. Discussion

We have identified a range of ML applications to sediment transport, morphologic, and hybrid coastal prediction problems. This provokes questions, such how to select an appropriate ML method to suit a particular problem, what sort of comparisons of ML can be made to more standard approaches, what are the common principles that can be applied to these applications, and where does this lead. These topics are discussed below.

6.1. Which ML technique and how much data?

We have reviewed commonly used ML techniques in the field of coastal morphodynamics, and investigated their applications across a range of scales. Fig. 5 is a schematic representation of our workflow. Some of the steps in this workflow are decisions (selecting an algorithm, determining if more data is needed, etc.) and may require significant investment of thought. With many possible tools, it is natural to ask if there is a specific ML routine that should be used for a given type of problem. There is no way of knowing this *a priori*, but trial and error informed by a desired outcome can guide ML routine selection. A researcher may have a preference for a given learning routine based on the type or quantity of data, an intuition regarding what worked best on a similar problem, or the form of a desired outcome. Guiding questions

can be used to determine the ML routine for a given problem: If probabilistic answers are required, researchers could focus on a probabilistic routine, such as Bayesian networks or many of the techniques that have gone unused by Coastal researchers (e.g., Gaussian processes, probabilistic interpretations of ANN, etc.). If multiple free parameters exist for a fixed, immutable equation, then a GA can be employed. If a specific smooth equation is needed (e.g., to be used in an analytical model), an ANN or GP can be used. If functional form and input dependencies are needed, a solution from GP can be attempted. Readers may note that there are trends in which ML routines are applied to a given problem — bars tend to be studied using ANN, while BN are often applied to shoreline and dune erosion. We do not believe that this occurs because one routine is better for solving a given task. Ambiguity over picking a ML approach for a given dataset highlights the need for a larger empirical study where many ML approaches are each attempted on an array of problems to determine (empirically) if there are optimal techniques for a given research question. Examples of this approach can be seen in other disciplines — work by Olson et al. (2017) and Hansen et al. (2013). Below (Section 6.3) we offer a set of best practices for coastal researchers who aspire to make their ML results usable for ML comparison.

It remains unclear how much data a researcher needs to perform meaningful ML analysis and produce a useful predictor. For example, Beuzen et al. (2017) investigated the amount of data needed to successfully parameterize a BN of shoreline change as a function of storm events, and suggested that the number likely depends on network complexity (degrees of freedom, or independent variables) and the signal clarity, likely a measure of signal to noise in the dataset. Goldstein and Coco (2014) and Tinoco et al. (2015) investigated the impact of adding data for GP prediction, finding that very little data was needed to train the model, and more data might actually provide degraded prediction for low-complexity models (Fig. 6). It is unclear if this relationship holds for other machine learning routines or prediction tasks, but could point the way in minimizing training data and maximizing the data used to test the resultant ML derived predictor. Similar to which algorithm, we expect that empirical guidance is most valuable, and further studies might shed light on how much data is really needed to make optimal predictions for a given problem, with a given prediction technique. Providing justification for a given amount of data, or a

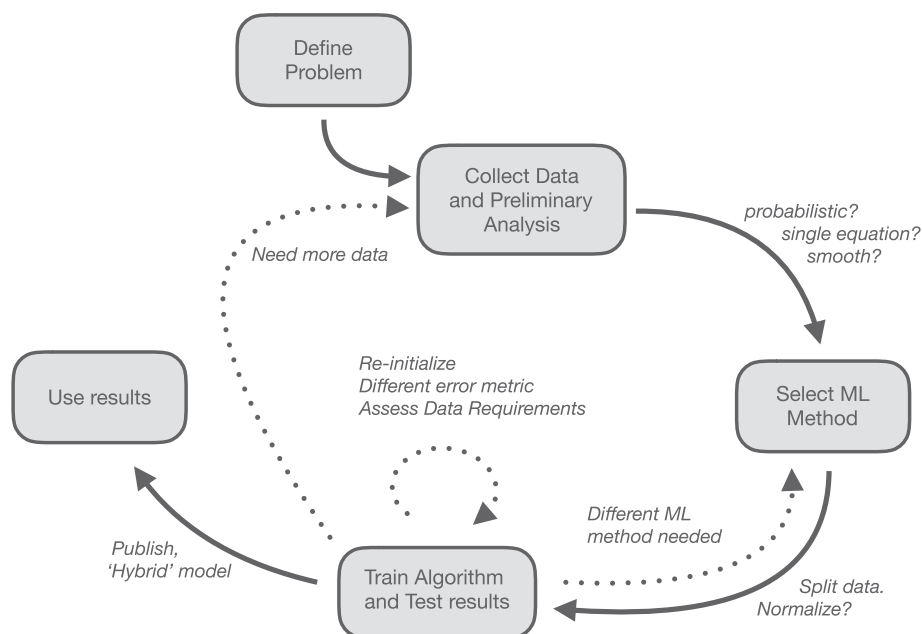


Fig. 5. A schematic of our workflow. Work starts with 'Define Problem' —new or published data — and progresses by following the filled arrows. Dotted arrows represent escapes from this workflow for various reasons.

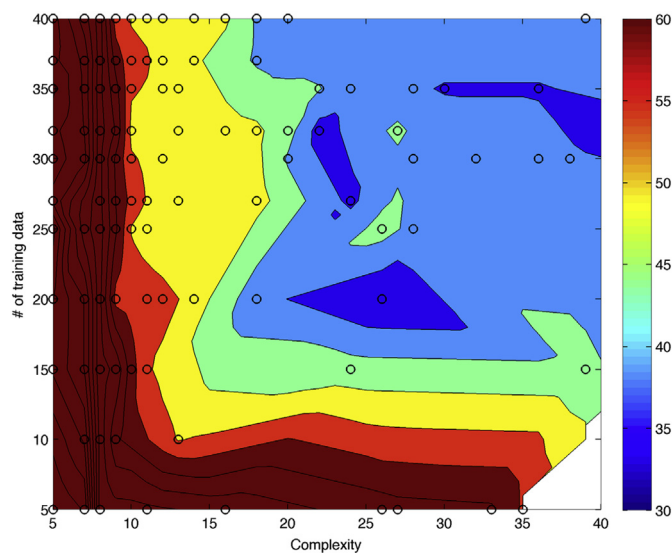


Fig. 6. A plot from Goldstein and Coco (2014) that displays GP solutions (open circles) of differing complexity developed by changing the size of the training data. The colormap is a measure of error, and suggests that continuing to increase the size of training data does not always yield decreases in error of the final solution.

given ML technique — especially quantitative justification — is particularly valuable for future researchers to determine what techniques and data quantities work for a given problem, and what can be learned (or what we are unable to learn) using a given approach and given set of data.

Beyond the amount of data needed for a ML analysis, all machine learning methods require that data be split discrete parts — with one part of the data used for developing the model (training data), while the remaining data is reserved outside the algorithm to test the developed model (testing data). The most important aspect of the data splitting is that users do not test the learned predictor with the training data (e.g., Domingos, 2012) — the testing data should not be seen by the ML algorithm as it is trained, testing must occur using a new, unseen portion of the data. There are enhancements beyond this simple two-part split (i.e., it is common to use additional subsets in the training process as validation to limit overfitting), but for this moment it is important that some portion of the data is used in the model building phase, and some portion of data is used in the model testing phase.

In many of the studies reviewed in this paper the data is split randomly. While random selection can be used, there are many methods to select a data for training and testing that are advantageous for developing generalizable predictor. For example, it is often easier to acquire data for coastal studies in controlled lab settings under weak forcing, and there may be few data points in field conditions (with complexities such as mixed grain sizes, irregular waves, etc.) or under extreme forcing (i.e., storms). Extracting training data by randomly sampling a full dataset may omit the critical information from the training data, for example extreme conditions (e.g., Passarella et al., 2018). There are many intentional sampling strategies from splitting data into training and testing that have not seen wide adoption in coastal community — we believe that adoption of these techniques will lead to more generalizable ML predictors. For examples of data selection routines, see Galelli et al. (2014) on evaluating input value selection routines, and previous work by Bowden et al. (2002), Bowden et al. (2012), Camus et al. (2011a,b), May et al. (2010), Tinoco et al. (2015), Splinter et al. (2013).

6.2. When does ML perform better than more traditional methods?

As we state in the sediment transport modeling section, the ML

studies we review often outperform more traditional prediction schemes that are based on derivations from first principles (i.e., conservation laws) and schemes based on more classic curve fitting techniques (e.g., linear regression). It remains unclear why this is the case. One possibility is that we are only aware of published ML studies, which bias us into believing ML provides only positive results. Or perhaps developing a new ML predictor might outperform a non-ML predictor because of flexibility — many ML routines can develop predictors that do not conform to a set basis function. For example, more traditional regression techniques work well when data conforms to a set functional form (a line, or a curve) and obeys the many generalizing assumptions (i.e., normally distributed). A researcher therefore performs dimensional reduction or transformation to first get data into this functional form, then fits the data with a known basis function (a line, a parabola, etc.). ML techniques offer more flexibility because the basis function for many techniques is highly adaptable — i.e., a neural network, which has many free parameters and can be trained to fit arbitrary curves. A further reason for ML outperforming traditional routines in the reviewed literature is that studies that take existing ML predictors and apply them to entirely new datasets, even those that are outside the range of data used to build the ML predictor, are rare. It remains to be seen how ML predictors will perform when data available for testing increase and more data from different sites and times are available for independent tests.

Many studies develop ML predictors for the same phenomena (e.g., suspended sediment transport flux, bar geometry, ripple wavelength). Even when multiple studies examine the same research question, the comparison between ML predictors is often impossible because each study uses different datasets, the method to split and testing data is often not clear, or information regarding the final predictor is not provided (e.g., for cases where ANN are used, the weights and biases may not be provided). Therefore it is difficult to truly compare ML methods and the resulting predictors. For instance, it is difficult to compare the ripple predictors developed by Yan and Zhang (2008) with a ANN and the ripple predictors developed by Goldstein et al. (2013) with a GP. After reviewing the work of ML in coastal morphodynamics and sediment transport we are left instead with the knowledge that the studies reviewed here can be understood as a proof of concept for ML being able to develop accurate predictions for a given dataset. ML is a clear way to develop bespoke prediction routines for a given site, a given dataset, or a given purpose. In some cases, more accurate prediction might be needed because of a specific research question or prediction task. An example is work by Goldstein et al. (2014), where a near bed suspended sediment reference concentration prediction scheme was developed from data collected in fields of sorted bedforms, and then used in a model of sorted bedforms to predict sorted bedform dynamics.

6.3. A set of practices for coastal ML research

If researchers using ML aspire to make their results reproducible, transferable, and useful to future researchers (e.g., intercomparison projects, helping to determine which algorithm is best, or how much data is needed), we offer a set of practices here that would aid in this goal. 1) Provide the data with the paper, or link to an open archive. 2) Unequivocally state the degree to which the training data and testing data are separate, that the testing data was not seen by the machine learning algorithm, and that the training data was not used to test the success of the developed model. 3) Clearly describe the technique used to split the data into training and testing, the percentage of data used for each group, and if possible, the actual data split into groups. If the data was additionally preprocessed (i.e., transformed, or some form of dimensionality reduction), explain this as well. 4) Report the final model in its entirety (e.g., weights, biases and architecture for ANN; binary splits for a tree model, etc.). 5) Define the metrics that are used to test the models and define levels required to be successful. 6)

Compare results to other models to provide benchmarks for improvements and the relative value of ML vs. theoretically developed models. 7) Compare ML model to newly collected data sets as a test to determine whether there are sufficient data for a particular model and whether the model is locally or generally applicable.

The goal with this set of best practices is to steer ML papers toward being usable and reproducible by the community. We understand that for a variety of reasons, these practices may not be possible under all circumstances, but if authors would like their ML work be built upon, tested, and refined, these practices aid in that goal. Very few studies that we review adhere to all of these practices. However we offer this guidance here to help advance the understanding and reuse of ML research in the coastal community, and we hope that all future work is performed with an eye toward reuse and reproducibility by others.

6.4. Future directions

Data-driven research relies on the existence of data. Beyond the collection and ad hoc sharing of data, the trend in publication of data is a major factor in the continued adoption of data-driven science. Data publication is enabled by the existence of repositories — Figshare, Zenodo, Data Dryad, Pangaea, and others identified by the re3data.org project (Pampel et al., 2013) — as well as data journals, publications that focus exclusively on descriptions of the data (from collection to access) such as Earth System Science Data (Pfeiffenberger and Carlson, 2011), Earth and Space Science (Hanson, 2014), Geoscience Data (Allan, 2014), and Scientific Data (Scientific Data, 2014). The continued collection and release of data will enable more data-driven work to occur in coastal settings.

In addition to data, ML research relies on ML algorithms and techniques. We have only discussed ML techniques that have been used in coastal settings, but this in no way is an exhaustive list of techniques. First, many common techniques — such as Support Vector Machines — have seen only minimal usage in coastal morphodynamics problems; though they have been used for coastal classification routines (Hoonhout et al., 2015) and in oceanographic contexts (Li et al., 2013). Second, many newer techniques may not yet have been applied for coastal research, such as recent advances in ANN architecture and training (e.g., Deep Learning; LeCun et al., 2015), or newer probabilistic techniques (Ghahramani, 2015). Third, automated machine learning (AutoML) has emerged as a method for algorithm selection and data preprocessing which could revolutionize the way in which Coastal scientists use ML. An example of AutoML is auto-sklearn — which is the work of Feurer et al. (2015) that is integrated with the Python package scikit-learn (Pedregosa et al., 2011). There is a world of new algorithms and techniques that can be brought over from the ML community — researchers might find it profitable to look for new techniques within the ML literature to make sure we are not missing out on the revolutionary advances in data-driven tools. In light of our general requests to continue publishing data for reuse, the continued adoption of new algorithms, and to teach these modern methods to students, we also see three specific areas for growth in Coastal ML research. First, to provide guidance for methods perform best for which problems, a more structured comparison projects between ML techniques using the identical data set (and data split) is required. Comparative work of this nature will help all researchers decide which ML has a high chance of success for a given problem. This may also help us to further understand a given coastal problem. Second, we identify an opportunity to use ML on timeseries, especially when systems may have memory and/or storage effects. The work of Pape et al. (2010) is a rare example, however new advances in neural network architecture (i.e., long short-term memory), and training has the potential to allow for more accurate time series prediction even when systems have strong autocorrelation, thresholds, and memory dynamic (e.g., shorelines, bars, bedforms, etc.) — see Kratzert et al. (2018) as an example of the power of long short-term memory networks in hydrological time series prediction. Third,

uncertainty derived from ML can further be incorporated into models. A clear possibility is to use probabilistic ML based predictions creatively in numerical models. An example is use the probabilistic nature of these predictors as stochastic parameterizations (e.g., Berner et al., 2017), whereby some aspect of an otherwise deterministic numerical model is made probabilistic, and models may then be able to generate ensemble predictions using identical forcing and initial conditions.

Finally, ML learning techniques can be thought of as new additions to the coastal researchers bag of tools — providing new insights during data analysis (e.g. Tinoco et al., 2015; Beuzen et al., 2018a). ML techniques could be taught alongside other more common data analysis techniques (e.g., Fourier transforms, wavelets; Zdeborová, 2017) since ultimately the goal is the same with any of these tools and techniques — to find and extract new knowledge or insight from data.

Acknowledgements

EBG gratefully acknowledges the support of UoA through a PBRF grant. GC funded by a GNS-Hazard Platform grant (contract 3710440). We are grateful for a review by Chris Sherwood at USGS, Daniel Buscombe, three anonymous reviewers, and comments by the ESR editor.

References

- Aguilera, P.A., Fernández, A., Fernández, R., Rumí, R., Salmerón, A., 2011. Bayesian networks in environmental modelling. *Environ. Model. Softw.* 26 (12), 1376–1388.
- Allan, R., 2014. Geoscience data. *Geosci. Data J.* 1 (1). <https://doi.org/10.1002/gdj3.3>.
- Alvarez-Ellacuria, A., Orfila, A., Gómez-Pujola, L., Simarro, G., Obregon, N., 2011. Decoupling spatial and temporal patterns in short-term beach shoreline response to wave climate. *Geomorphology* 128, 199–208.
- Anderson, C., 2008. The end of theory: the data deluge makes the scientific method obsolete. *Wired Mag.* 16 (7), 07–16.
- Apotos, A., Raubenheimer, B., Elgar, S., Guza, R.T., 2008. Testing and calibrating parametric wave transformation models on natural beaches. *Coast. Eng.* 55, 224–235.
- Aubrey, D.G., 1979. Seasonal patterns of onshore/offshore sediment movement. *J. Geophys. Res.* 84, 6347–6354.
- Babovic, V., Keijzer, M., 2000. Genetic programming as a model induction engine. *J. Hydroinf.* 2 (1), 35–60.
- Berner, J., Achatz, U., Batte, L., Bengtsson, L., Cámara, A.D.L., Christensen, H.M., ... Franzke, C.L., 2017. Stochastic parameterization: toward a new view of weather and climate models. *Bull. Am. Meteorol. Soc.* 98 (3), 565–588.
- Beuzen, T., Simmons, J., 2019. A variable selection package driving Netica with Python. *Environ. Model. Softw.* 115, 1–5. <https://doi.org/10.1016/j.envsoft.2019.01.018>.
- Beuzen, T., Splinter, K.D., Turner, I.L., Harley, M.D., Marshall, L., 2017. Predicting storm erosion on sandy coastlines using a Bayesian network. *Australas. Coasts Ports 2017* (Working with Nature, 102).
- Beuzen, T., Splinter, K.D., Marshall, L.A., Turner, I.L., Harley, M.D., Palmsten, M.L., 2018a. Bayesian Networks in coastal engineering: distinguishing descriptive and predictive applications. *Coast. Eng.* 135, 16–30.
- Beuzen, T., Marshall, L., Splinter, K.D., 2018b. A comparison of methods for discretizing continuous variables in Bayesian Networks. *Environ. Model. Softw.* 108, 61–66. <https://doi.org/10.1016/j.envsoft.2018.07.007>.
- Bilskie, M.V., Hagen, S.C., Alizad, K., Medeiros, S.C., Passeri, D.L., Needham, H.F., Cox, A., 2016. Dynamic simulation and numerical analysis of hurricane storm surge under sea level rise with geomorphologic changes along the northern Gulf of Mexico. *Earth's Future* 4, 177–193. <https://doi.org/10.1002/2015EF000347>.
- Bolaños, R., Souza, A., 2010. Measuring hydrodynamics and sediment transport processes in the Dee Estuary. *Earth Syst. Sci. Data* 2, 157–165. <https://doi.org/10.5194/essd-2-157-2010>.
- Bowden, G.J., Maier, H.R., Dandy, G.C., 2002. Optimal division of data for neural network models in water resources applications. *Water Resour. Res.* 38 (2).
- Bowden, G.J., Maier, H.R., Dandy, G.C., 2012. Real-time deployment of artificial neural network forecasting models: Understanding the range of applicability. *Water Resour. Res.* 48 (10).
- Bulteau, T., Baills, A., Petitjean, L., Garcin, M., Palanisamy, H., Le Cozannet, G., 2015. Gaining insight into regional coastal changes on La Réunion island through a Bayesian data mining approach. *Geomorphology* 228, 134–146.
- Camus, P., Mendez, F.J., Medina, R., Cofiño, A.S., 2011a. Analysis of clustering and selection algorithms for the study of multivariate wave climate. *Coast. Eng.* 58 (6), 453–462.
- Camus, P., Cofiño, A.S., Mendez, F.J., Medina, R., 2011b. Multivariate wave climate using self-organizing maps. *J. Atmos. Ocean. Technol.* 28 (11), 1554–1568.
- Charniak, E., 1991. Bayesian networks without tears. *AI Mag.* 12 (4), 50–63.
- Chau, K., 2006. A review on the integration of artificial intelligence into coastal modeling. *J. Environ. Manag.* 80 (1), 47–57.
- Chen, S.H., Pollino, C.A., 2012. Good practice in Bayesian network modelling. *Environ. Model. Softw.* 37, 134–145.

- Cohen, I.B., Whitman, A., Budenz, J., 2016. The Principia: The Authoritative Translation and Guide: Mathematical Principles of Natural Philosophy.
- Cooper, G.E., Herskovits, E., 1992. A Bayesian method for the induction of probabilistic networks from data. *Mach. Learn.* 9, 309–347.
- Coveney, P.V., Dougherty, E.R., Highfield, R.R., 2016. Big data need big theory too. *Phil. Trans. R. Soc. A* 374 (2080), 20160153.
- Crutchfield, J.P., 2014. The dreams of theory. *Wiley Interdiscip. Rev.* 6 (2), 75–79.
- D'Ambrosio, D., Spataro, W., Rongo, R., Iovine, G., 2013. Genetic algorithms, optimization, and evolutionary modeling. In: Shroder, J. (Editor in Chief), Baas, A. C. W. (Ed.), *Quantitative Modeling of Geomorphology*. Academic Press, San Diego, CA, vol. vol. 2, pp. 74–97.
- Data, Scientific, 2014. More bang for your byte. *Sci. Data* 1, 140010.
- De'ath, G., 2007. Boosted trees for ecological modeling and prediction. *Ecology* 88 (1), 243–251.
- De'ath, G., Fabricius, K.E., 2000. Classification and regression trees: a powerful yet simple technique for ecological data analysis. *Ecology* 81 (11), 3178–3192.
- Demirci, M., Üneş, F., Aköz, M.S., 2015. Prediction of cross-shore sandbar volumes using neural network approach. *J. Mar. Sci. Technol.* 20 (1), 171–179.
- den Heijer, C.K., Knipping, D.T., Plant, N.G., de Vries, J.S.V.T., Baart, F., van Gelder, P.H., 2012. Impact assessment of extreme storm events using a Bayesian network. *Coast. Eng. Proc.* 1 (33), 4.
- Dickson, M.E., Perry, G.L., 2016. Identifying the controls on coastal cliff landslides using machine-learning approaches. *Environ. Model. Softw.* 76, 117–127.
- Domingos, P., 2012. A few useful things to know about machine learning. *Commun. ACM* 55 (10), 78–87.
- Elith, J., Leathwick, J.R., Hastie, T., 2008. A working guide to boosted regression trees. *J. Anim. Ecol.* 77, 802–813.
- Feurer, M., Klein, A., Eggenberger, K., Springenberg, J., Blum, M., Hutter, F., 2015. Efficient and robust automated machine learning. In: *Advances in neural information processing systems* 28, pp. 2962–2970. <http://papers.nips.cc/paper/5872-efficient-and-robust-automated-machine-learning.pdf>.
- Galelli, S., Humphrey, G.B., Maier, H.R., Castelletti, A., Dandy, G.C., Gibbs, M.S., 2014. An evaluation framework for input variable selection algorithms for environmental data-driven models. *Environ. Model. Softw.* 62, 33–51.
- Garel, E., Ferreira, Ó., 2015. Multi-year high-frequency physical and environmental observations at the Guadiana Estuary. *Earth Syst. Sci. Data* 7, 299–309. <https://doi.org/10.5194/essd-7-299-2015>.
- Ghahramani, Z., 2015. Probabilistic machine learning and artificial intelligence. *Nature* 521 (7553), 452–459.
- Goldstein, E.B., Coco, G., 2014. A machine learning approach for the prediction of settling velocity. *Water Resour. Res.* 50 (4), 3595–3601.
- Goldstein, E.B., Coco, G., 2015. Machine learning components in deterministic models: hybrid synergy in the age of data. *Front. Environ. Sci.* 3, 33.
- Goldstein, E.B., Moore, L.J., 2016. Stability and bistability in a one-dimensional model of coastal foredune height. *J. Geophys. Res.* 121 (5), 964–977.
- Goldstein, E.B., Moore, L.J., 2018. A calibration workflow for coastal dune models. *Shore Beach* 86 (3), 47–51. available at. [10.31223/osf.io/cd87u](https://doi.org/10.31223/osf.io/cd87u).
- Goldstein, E.B., Coco, G., Murray, A.B., 2013. Prediction of wave ripple characteristics using genetic programming. *Cont. Shelf Res.* 71, 1–15. <https://doi.org/10.1016/j.csr.2013.09.020>.
- Goldstein, E.B., Coco, G., Murray, A.B., Green, M.O., 2014. Data driven components in a model of inner shelf sorted bedforms: a new hybrid model. *Earth Surf. Dynam.* Discuss. 1, 531–569. <https://doi.org/10.5194/esurf-d-1-531-2013>.
- Grimes, D.J., Cortale, N., Baker, K., McNamara, D.E., 2015. Nonlinear forecasting of intertidal shelf evolution. *Chaos* 25 (10), 103116.
- Gutiérrez, B.T., Plant, N.G., Thieler, E.R., 2011. A Bayesian network to predict the coastal vulnerability to sea-level rise. *J. Geophys. Res.* 116, F02009. <https://doi.org/10.1029/2010JF001891>.
- Gutiérrez, B.T., Plant, N.G., Pendleton, E.A., Thieler, E.R., 2014. Using a Bayesian network to predict shoreline change vulnerability to sea-level rise for the coasts of the United States. *U.S. Geol. Surv. Open File Rep.* (2014–108326pp., U.S. Geological Survey, Reston, Va).
- Gutiérrez, B.T., Plant, N.G., Thieler, E.R., Turecek, A., 2015. Using a Bayesian network to predict barrier island geomorphologic characteristics. *J. Geophys. Res.* 120 (12), 2452–2475.
- Hansen, K., Montavon, G., Biegler, F., Fazli, S., Rupp, M., Scheffler, M., Müller, K.R., 2013. Assessment and validation of machine learning methods for predicting molecular atomization energies. *J. Chem. Theory Comput.* 9 (8), 3404–3419.
- Hanson, B., 2014. AGU to launch a new open-access journal spanning the earth and space sciences. *EOS Trans. Am. Geophys. Union* 95 (6), 56.
- Hapke, C., Plant, N.G., 2010. Predicting coastal cliff erosion using a bayesian probabilistic model. *Mar. Geol.* 278, 140–149. <https://doi.org/10.1016/j.margeo.2010.10.001>.
- Hashemi, M.R., Ghadampour, Z., Neill, S.P., 2010. Using an artificial neural network to model seasonal changes in beach profiles. *Ocean Eng.* 37, 1345–1356.
- Hastie, T., Tibshirani, R., Friedman, J., 2009. The elements of statistical learning: data mining, inference, and prediction. (Springer series in statistics).
- Hey, T., Tansley, S., Tolle, K.M., 2009. The Fourth Paradigm: data-intensive Scientific Discovery. Microsoft research, Redmond, WA.
- Holland, J.H., 1975. Adaptation in natural and artificial systems. An introductory analysis with application to biology, control, and artificial intelligence. Ann Arbor, MI. University of Michigan Press.
- Hoonhout, B.M., Rademacher, M., Baart, F., Van der Maaten, L.J.P., 2015. An automated method for semantic classification of regions in coastal images. *Coast. Eng.* 105, 1–12.
- Hsieh, W.W., 2004. Nonlinear multivariate and time series analysis by neural network methods. *Rev. Geophys.* 42, RG1003. <https://doi.org/10.1029/2002RG000112>.
- Hsieh, W.W., 2009. Machine learning methods in the environmental sciences: Neural networks and kernels. Cambridge university press.
- Iglesias, G., Lopez, I., Carballo, R., Castro, A., 2009a. Headland-bay beach planform and tidal range: a neural network model. *Geomorphology* 112, 135–143.
- Iglesias, G., López, I., Castro, A., Carballo, R., 2009b. Neural network modelling of planform geometry of headland-bay beaches. *Geomorphology* 103 (4), 577–587.
- Iglesias, G., Diz-Lois, G., Pinto, F.T., 2010. Artificial Intelligence and headland-bay beaches. *Coast. Eng.* 57 (2), 176–183.
- Jaffe, B.E., Rubin, D.M., 1996. Using nonlinear forecasting to learn the magnitude and phasing of time-varying sediment suspension in the surf zone. *J. Geophys. Res.* 101 (C6), 14,283–14,296.
- Jäger, W.S., Christie, E.K., Hanea, A.M., den Heijer, C., Spencer, T., 2018. A Bayesian network approach for coastal risk analysis and decision making. *Coast. Eng.* 134, 48–61.
- Jones, N.S., MacCarone, T.J., 2013. Inference for the physical sciences. *Phil. Trans. R. Soc. A* 371, 20120493. <https://doi.org/10.1098/rsta.2012.0493>.
- Kabiri-Samani, A.R., Aghaee-Tarazjani, J., Borghai, S.M., Jeng, D.S., 2011. Application of neural networks and fuzzy logic models to long-shore sediment transport. *Appl. Soft Comput.* 11, 2880–2887.
- Kingston, K.S., Ruessink, B.G., van Enckevort, I.M.J., Davidson, M.A., 2000. Artificial neural network correction of remotely sensed sandbar location. *Mar. Geol.* 169, 137–160.
- Kizhisseri, A.S., Simmonds, D., Rafiq, Y., Borthwick, M., 2005. An evolutionary computation approach to sediment transport modeling. In: *Fifth International Conference on Coastal Dynamics*. Spain, Barcelona April 4–8, 2005.
- Knaepen, M.A.F., Hulscher, S.J.M.H., 2002. Regeneration of sand waves after dredging. *Coast. Eng.* 46, 277–289.
- Komurcu, M.I., Tutkun, N., Ozolcer, I.H., Akpinar, A., 2008. Estimation of the beach bar parameters using the genetic algorithms. *Appl. Math. Comput.* 195, 49–60.
- Komurcu, M.I., Komur, M.A., Akpinar, A., Ozolcer, I.H., Yuksek, O., 2013. Prediction of offshore bar-shape parameters resulting by cross-shore sediment transport using artificial neural network. *Appl. Ocean Res.* 40, 74–82.
- Koza, J.R., 1992. Genetic programming, on the programming of computers by means of natural selection. In: MIT Press. USA, Cambridge, MA.
- Krasnopolsky, V.M., 2013. The Application of Neural Networks in the Earth Sciences. Springer.
- Krasnopolsky, V.M., Fox-Rabinovitz, M.S., 2006. A new synergetic paradigm in environmental numerical modeling: Hybrid models combining deterministic and machine learning components. *Ecol. Model.* 191, 5–18.
- Kratzert, F., Klotz, D., Brenner, C., Schulz, K., Herrnegger, M., 2018. Rainfall-Runoff modelling using Long-Short-Term-Memory (LSTM) networks. *Hydrol. Earth Syst. Sci. Discuss.* <https://doi.org/10.5194/hess-2018-247>.
- Lazarus, E.D., McNamara, D.E., Smith, M.D., Gopalakrishnan, S., Murray, A.B., 2011. Emergent behavior in a coupled economic and coastline model for beach nourishment. *Nonlinear Process. Geophys.* 18 (6), 989–999.
- LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *Nature* 521 (7553), 436–444.
- Lentz, E.E., Hapke, C.J., 2011. Geologic framework influences on the geomorphology of an anthropogenically modified barrier island: assessment of dune/beach changes at Fire Island, New York. *Geomorphology* 126 (1), 82–96.
- Lentz, E.E., Thieler, E.R., Plant, N.G., Stippa, S.R., Horton, R.M., Gesch, D.B., 2016. Evaluation of dynamic coastal response to sea-level rise modifies inundation likelihood. *Nat. Clim. Chang.* 6 (7), 696–700.
- Li, Z., Li, L., Song, K., Cassar, N., 2013. Estimation of phytoplankton size fractions based on spectral features of remote sensing ocean color data. *J. Geophys. Res. Oceans* 118. <https://doi.org/10.1002/jgrc.20137>.
- Limber, P.W., Murray, A.B., 2014. Unraveling the dynamics that scale cross-shore headland relief on rocky coastlines: 2. Model predictions and initial tests. *J. Geophys. Res.* 119 (4), 874–891.
- Limber, P.W., Brad Murray, A., Adams, P.N., Goldstein, E.B., 2014. Unraveling the dynamics that scale cross-shore headland relief on rocky coastlines: 1. Model development. *J. Geophys. Res.* 119 (4), 854–873.
- Lin, S., Sheng, J., 2017. Assessing the performance of wave breaking parameterizations in shallow waters in spectral wave models. *Ocean Model.* 120, 41–59.
- López, I., Aragonés, L., Villacampa, Y., Serra, J.C., 2017. Neural network for determining the characteristic points of the bars. *Ocean Eng.* 136, 141–151.
- Lorenz, E.N., 1969a. Atmospheric predictability as revealed by naturally occurring analogues. *J. Atmos. Sci.* 26, 636–646.
- Lorenz, E.N., 1969b. Three approaches to atmospheric predictability. *Bull. Am. Meteorol.* 50, 345–349.
- Loureiro, C., Ferrera, O., Cooper, J.A., 2013. Applicability of parametric beach morphodynamic state classification on embayed beaches. *Mar. Geol.* 346, 153–164.
- Luijendijk, A., Hagenaars, G., Ranasinghe, R., Baart, F., Donchyts, G., Aarninkhof, S., 2018. The State of the World's Beaches. *Sci. Rep.* 8 (6641). <https://doi.org/10.1038/s41598-018-24630-6>.
- Mafi, S., Yeganeh-Bakhtiary, A., Kazeminezhad, M.H., 2013. Prediction formula for longshore sediment transport rate with M5' algorithm. In: Mafi, S., Yeganeh-Bakhtiary, A., Kazeminezhad, M.H. (Eds.), *Proceedings 12th International Coastal Symposium (Plymouth, England)*, Journal of Coastal Research, pp. 2149–2154 Special Issue No. 65. (ISSN 0749-0208).
- Malakoff, D., 1999. Bayes offers a 'New' way to make sense of numbers. *Science* 286 (5444), 1460–1464.
- May, R.J., Maier, H.R., Dandy, G.C., 2010. Data splitting for artificial neural networks using SOM-based stratified sampling. *Neural Netw.* 23 (2), 283–294.
- Mitchell, M., 1995. Genetic algorithms: An overview. *Complexity* 1 (1), 31–39.
- Mitchell, M., 1998. An introduction to genetic algorithms. MIT press.
- Múnera, S., Osorio, A.F., Velásquez, J.D., 2014. Data-based methods and algorithms for

- the analysis of sandbar behavior with exogenous variables. *Comput. Geosci.* 72, 134–146.
- Murray, A.B., Lazarus, E., Ashton, A., Baas, A., Coco, G., Coulthard, T., Fonstad, M., Haff, P., McNamara, D., Paola, C., Pelletier, J., Reinhardt, L., 2009. Geomorphology, complexity, and the emerging science of the Earth's surface. *Geomorphology* 103 (3), 496–505.
- Murray, A.B., Goldstein, E.B., Coco, G., 2014a. Cause and effect in geomorphic systems: complex-systems perspectives. *Geomorphology* 219, 1–9.
- Murray, A.B., Coco, G., Goldstein, E.B., 2014b. The shape of patterns to come: from initial formation to long-term evolution. *Earth Surf. Process. Landf.* <https://doi.org/10.1002/esp.3487>.
- Murray, A.B., Gasparini, N.M., Goldstein, E.B., van der Wegen, M., 2016. Uncertainty quantification in modeling earth surface processes: more applicable for some types of models than for others. *Comput. Geosci.* 90 (Part B), 6–16. <https://doi.org/10.1016/j.cageo.2016.02.008>.
- Nelson, T.R., Voulgaris, G., Traykovski, P., 2013. Predicting wave-induced ripple equilibrium geometry. *J. Geophys. Res.* 118 (6), 3202–3220.
- Nylén, T., Hellemaa, P., Luoto, M., 2015. Determinants of sediment properties and organic matter in beach and dune environments based on boosted regression trees. *Earth Surf. Process. Landf.* 40 (9), 1137–1145.
- Oehler, F., Coco, G., Green, M.O., Bryan, K.R., 2011. A data-driven approach to predict suspended-sediment reference concentration under non-breaking waves. *Cont. Shelf Res.* 46, 96–106.
- Olden, J.D., Joy, M.K., Death, R.G., 2004. An accurate comparison of methods for quantifying variable importance in artificial neural networks using simulated data. *Ecol. Model.* 178, 389–397.
- Olden, J.D., Lawler, J.J., Poff, N.L., 2008. Machine learning methods without tears: a primer for ecologists. *Q. Rev. Biol.* 83 (2), 171–193.
- Olson, R.S., La Cava, W., Mustahsan, Z., Varik, A., Moore, J.H., 2017. Data-driven advice for applying machine learning to bioinformatics problems. *arXiv preprint arXiv:1708.05070*.
- O'Neill, M., Vanneschi, L., Gustafson, S., Banzhaf, W., 2010. Open issues in genetic programming. *Genet. Program. Evol. M.* 11, 339–363.
- Palmsten, M.L., Splinter, K.D., Plant, N.G., Stockdon, H.F., 2014. Probabilistic estimation of dune retreat on the Gold Coast, Australia. *Shore Beach* 82 (4), 35–43.
- Pampel, H., Vierkant, P., Scholze, F., Bertelmann, R., Kindling, M., Klump, J., Dierolf, U., 2013. Making research data repositories visible: the re3data.org registry. *PLoS ONE* 8 (11), e78080.
- Pape, L., Ruessink, B.G., Wiering, M.A., Turner, I.L., 2007. Recurrent neural network modeling of nearshore sandbar behavior. *Neural Netw.* 20, 509–518.
- Pape, L., Kuriyama, Y., Ruessink, B.G., 2010. Models and scales for cross-shore sandbar migration. *J. Geophys. Res.* 115, F03043. <https://doi.org/10.1029/2009JF001644>.
- Passarella, M., Goldstein, E.B., De Muro, S., Coco, G., 2018. The use of genetic programming to develop a predictor of swash excursion on sandy beaches. *Nat. Hazards Earth Syst. Sci.* 18, 599–611. <https://doi.org/10.5194/nhess-18-599-2018>.
- Passeri, D.L., Hagen, S.C., Plant, N.G., Bilske, M.V., Medeiros, S.C., Alizad, K., 2016. Tidal hydrodynamics under future sea level rise and coastal morphology in the Northern Gulf of Mexico. *Earth's Future* 4, 159–176. <https://doi.org/10.1002/2015EF000332>.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al., 2011. Scikit-learn: machine learning in python. *J. Mach. Learn. Res.* 12, 2825–2830.
- Pfeifferberger, H., Carlson, D., 2011. "Earth System Science Data" (ESSD)-A Peer Reviewed Journal for Publication of Data. *D-Lib Mag.* 17 (1/2).
- Pinsky, M.L., Guannel, G., Arkema, K.K., 2013. Quantifying wave attenuation to inform coastal habitat conservation. *Ecosphere* 4 (8).
- Plant, N.G., Holland, K.T., 2011. Prediction and assimilation of surf-zone processes using a Bayesian network. Part I: Forward models. *Coast. Eng.* 58 (1), 119–130.
- Plant, N.G., Stockdon, H.F., 2012. Probabilistic prediction of barrier-island response to hurricanes. *J. Geophys. Res.* 117, F03015. <https://doi.org/10.1029/2011JF002326>.
- Plant, N.G., Holman, R.A., Freilich, M.H., Birkemeier, W.A., 1999. A simple model for interannual sandbar behavior. *J. Geophys. Res.* 104 (C7), 15755–15776.
- Plant, N.G., Holland, K.T., Holman, R.A., 2006. A dynamical attractor governs beach response to storms. *Geophys. Res. Lett.* 33 (17).
- Plant, N.G., Aarninkhof, S.G.J., Turner, I.L., Kingston, K., 2007. The performance of shoreline detection models applied to video imagery. *J. Coast. Res.* 23 (3), 658–670.
- Plant, N.G., Flocks, J., Stockdon, H.F., Long, J.W., Guy, K., Thompson, D.M., Cormier, J.M., Smith, C.G., Miselis, J.L., Dalyander, P.S., 2014. Predictions of barrier island berm evolution in a time-varying storm climatology. *J. Geophys. Res. Earth Surf.* 119, 300–316. <https://doi.org/10.1002/2013JF002871>.
- Plant, N.G., Thieler, E.R., Passeri, D.L., 2016. Coupling centennial-scale shoreline change to sea-level rise and coastal morphology in the Gulf of Mexico using a Bayesian network. *Earth's Future* 4, 1. <https://doi.org/10.1002/2015EF000331>.
- Plomaritis, T.A., Costas, S., Ferreira, O., 2018. Use of a Bayesian Network for coastal hazards, impact and disaster risk reduction assessment at a coastal barrier (Ria Formosa, Portugal). *Coast. Eng.* 134, 134–147.
- Poelhekke, L., Jäger, W.S., van Dongeren, A., Plomaritis, T.A., McCall, R., Ferreira, O., 2016. Predicting coastal hazards for sandy coasts with a Bayesian network. *Coast. Eng.* 118, 21–34.
- Poli, R., Langdon, W.B., McPhee, N.F., 2008. A field guide to genetic programming. (Lulu Enterprises UK Limited).
- Popper, K., Miller, D., 1983. A proof of the impossibility of inductive probability. *Nature* 302, 687–688.
- Rigos, A., Tsekouras, G.E., Voudoukas, M.I., Chatzipavlis, A., Velegrakis, A.F., 2016a. A Chebyshev polynomial radial basis function neural network for automated shoreline extraction from coastal imagery. *Integrated Computer-Aided Engineering* 23 (2), 141–160.
- Rigos, A., Tsekouras, G.E., Chatzipavlis, A., Velegrakis, A.F., 2016b. Modeling beach rotation using a novel legendre polynomial feedforward neural network trained by nonlinear constrained optimization. In: *IFIP International Conference on Artificial Intelligence Applications and Innovations*. Springer International Publishing, pp. 167–179.
- Roelvink, D., Reniers, A., Van Dongeren, A.P., de Vries, J.V.T., McCall, R., Lescinski, J., 2009. Modelling storm impacts on beaches, dunes and barrier islands. *Coast. Eng.* 56 (11), 1133–1152.
- Rubin, David M., 1992. Use of forecasting signatures to help distinguish periodicity, randomness, and chaos in ripples and other spatial patterns. *Chaos* 2 (4), 525–535.
- Ruessink, B.G., 2005. Calibration of nearshore process models: application of a hybrid genetic algorithm. *J. Hydroinf.* 7, 135–149.
- Schmidt, M., Lipson, H., 2009. Distilling free-form natural laws from experimental data. *Science* 324, 81–85.
- Sherman, D.J., 1995. Problems of scale in the modeling and interpretation of coastal dunes. *Mar. Geol.* 124 (1–4), 339–349.
- Splinter, K.D., Holman, R.A., Plant, N.G., 2011. A behavior-oriented dynamic model for sandbar migration and 2DH evolution. *J. Geophys. Res.* C 116 (1).
- Splinter, K.D., Turner, I.L., Davidson, M.A., 2013. How much data is enough? The importance of morphological sampling interval and duration for calibration of empirical shoreline models. *Coast. Eng.* 77, 14–27.
- Stalzer, M., Mentzel, C., 2016. A preliminary review of influential works in data-driven discovery. *SpringerPlus* 5 (1), 1266.
- Stephens, S.A., Coco, G., Bryan, K.R., 2011. Numerical simulations of wave setup over barred beach profiles: implications for predictability. *J. Waterw. Port Coast. Ocean Eng.* 137 (4), 175–181.
- Stockdon, H.F., Thompson, D.M., Plant, N.G., Long, J.W., 2014. Evaluation of wave runup predictions from numerical and parametric models. *Coast. Eng.* 92, 1–11.
- Tinoco, R.O., Goldstein, E.B., Coco, G., 2015. A data-driven approach to develop physically sound predictors: Application to depth-averaged velocities on flows through submerged arrays of rigid cylinders. *Water Resour. Res.* 51 (2), 1247–1263.
- Tsekouras, G.E., Rigos, A., Chatzipavlis, A., Velegrakis, A., 2015. A neural-fuzzy network based on Hermite polynomials to predict the coastal erosion. In: *Engineering Applications of Neural Networks*. Springer International Publishing, pp. 195–205.
- Turner, I.L., Harley, M.D., Short, A.D., Simmons, J.A., Bracs, M.A., Phillips, M.S., Splinter, K.D., 2016. A multi-decade dataset of monthly beach profile surveys and inshore wave forcing at Narrabeen, Australia. *Scientific Data* 3.
- Uusitalo, L., 2007. Advantages and challenges of Bayesian networks in environmental modelling. *Ecol. Model.* 203 (3), 312–318.
- Valentine, A.P., Kalnins, L.M., 2016. An introduction to learning algorithms and potential applications in geomorphometry and earth surface dynamics. *Earth surface dynamics*. 4, 445–460.
- van der Werf, J.J., Schretlen, J.J., Ribberink, J.S., O'Donoghue, T., 2009. Database of full-scale laboratory experiments on wave-driven sand transport processes. *Coast. Eng.* 56 (7), 726–732.
- van Maanen, B., Coco, G., Bryan, K.R., Ruessink, B.G., 2010. The use of artificial neural networks to analyze and predict alongshore sediment transport. *Nonlinear Process. Geophys.* 17, 395–404. <https://doi.org/10.5194/npg-17-395-2010>.
- Van Verseveld, H.C.W., Van Dongeren, A.R., Plant, N.G., Jäger, W.S., den Heijer, C., 2015. Modelling multi-hazard hurricane damages on an urbanized coast with a Bayesian Network approach. *Coast. Eng.* 103, 1–14.
- Werner, B.T., 1999. Complexity in natural landform patterns. *Science* 284 (5411), 102–104.
- Wikle, C.K., Berliner, L.M., 2007. A Bayesian tutorial for data assimilation. *Physica D* 230, 1–16.
- Wilson, K.E., Adams, P.N., Hapke, C.J., Lentz, E.E., Brenner, O., 2015. Application of Bayesian Networks to hindcast barrier island morphodynamics. *Coast. Eng.* 102, 30–43.
- Winant, C.D., Inman, D.L., Nordstrom, C.E., 1975. Description of seasonal beach changes using empirical eigenfunctions. *J. Geophys. Res.* 80 (15), 1979–1986.
- Wright, L.D., Short, A.D., 1984. Morphodynamic variability of surf zones and beaches: a synthesis. *Mar. Geol.* 56 (1–4), 93–118.
- Wright, L.D., Short, A.D., Green, M.O., 1985. Short-term changes in the morphodynamic states of beaches and surf zones: An empirical predictive model. *Mar. Geol.* 62, 339–364.
- Yan, B., Zhang, Q., Wai, O.W.H., 2008. Prediction of sand ripple geometry under waves using an artificial neural network. *Comput. Geosci.* 34, 1655–1664.
- Yao, X., 1999. Evolving artificial neural networks. *Proc. IEEE* 87 (9), 1423–1447.
- Yates, M.L., Le Cozannet, G., 2012. Evaluating European coastal evolution using Bayesian networks. *Nat. Hazards Earth Syst. Sci.* 12, 1173–1177.
- Yates, M.L., Guza, R.T., O'Reilly, W.C., 2009. Equilibrium shoreline response: Observations and modeling. *J. Geophys. Res.* C 114 (9).
- Yoon, H.-D., Cox, D.T., Kim, M., 2013. Prediction of time-dependent sediment suspension in the surf zone using artificial neural network. *Coast. Eng.* 71, 78–86. <https://doi.org/10.1016/j.coastaleng.2012.08.005>.
- Zdeborová, L., 2017. Machine learning: new tool in the box. *Nat. Phys.* 13 (5), 420–421.