1      **Title Page**
2
3      **Extracting Shallow-water Bathymetry from Lidar Point Clouds Using Pulse Attribute Data: Merging**
4      **Density-based and Machine Learning Approaches**
5

6      Kim Lowell[1], Brian Calder
7      Center for Coastal and Ocean Mapping and Joint Hydrographic Center
8      University of New Hampshire
9      24 Colovos Road,
10     Durham, NH 03824   UNITED STATES
11     [1]**Corresponding Author**
12
13     **Titles and email addresses**:
14         Kim Lowell, Research Data Scientist, klowell@ccom.unh.edu (ORCID: 000-0002-8326-4022)
15         Brian Calder, Research Professor, brc@ccom.unh.edu
16

19     **Disclosure**: The authors have no potential competing financial nor non-financial interests in the work

20     presented.

21     **Data availability statement**: Data that support the findings of this study are available at the link
22     doi.org/10.6084/m9.figshare.12597404.  SBET data in the required format are provided at the figshare
23     link.  Though the .las data used are available to the public, the authors are not authorized to make them
24     directly available.  A small sample of the data for a single data tile are provided at the figshare link.
25     Complete      data      sets      (2016_420500e_2728500n.laz,      2016_426000e_2708000n.laz,
26     2016_428000e_2719500n.laz,   and   2016_430000e_2707500n.laz)   can   be   downloaded   from
27     https://coast.noaa.gov/digitalcoast/data/ (Data set name: 2016 NGS Topobathy Lidar: Key West FL') as
28     compressed .laz files.  These can be decompressed using the LASzip tool which can be downloaded from
29     laszip.org.
30
31     **Total word count**: 9600          **Word count excluding title page and references**: 7746
32

33 **Extracting Shallow-water Bathymetry from Lidar Point Clouds Using Pulse Attribute Data: Merging**
34 **Density-based and Machine Learning Approaches**
35
36 **Abstract**

37 To automate extraction of bathymetric soundings from lidar point clouds, two machine learning (ML[1])
38 techniques were combined with a more conventional density-based algorithm.  The study area was four
39 data "tiles" near the Florida Keys.  The density-based algorithm determined the most likely depth (MLD)
40 for a grid of "estimation nodes" (ENs).  Unsupervised $k$-means clustering determined which EN's MLD
41 depth and associated soundings represented ocean depth rather than ocean surface or noise to produce
42 a preliminary classification.  An extreme gradient boosting (XGB) model was fitted to pulse return
43 metadata – e.g., return intensity, incidence angle -- to produce a final *Bathy/NotBathy* classification.
44 Compared to an operationally produced reference classification, the XGB model increased global accuracy
45 and decreased the false negative rate (FNR) – i.e., undetected bathymetry – that are most important for
46 nautical navigation for all but one tile.  Agreement between the final XGB and operational reference
47 classifications ranged from 0.84 to 0.999.  Imbalance between *Bathy* and *NotBathy* was addressed using
48 a probability decision threshold that equalizes the FNR and the true positive rate (TPR).  Two methods are
49 presented for visually evaluating differences between the two classifications spatially and in feature-
50 space.

51
52 **Keywords**: shallow water bathymetry, airborne lidar, Florida Keys, extreme gradient boosting, k-means
53 clustering
54

55 **1.   Introduction and Approach**

56 It is generally accepted that Hickman and Hogg (1969) authored the first article published on the use of
57 airborne lidar ('light detection and ranging') data for bathymetric mapping.  They observed that due to
58 limitations on the penetration of light through water, lidar is most appropriate for shallow water charting.
59 Heritage and Hetherington (2007) noted that the initial focus of lidar research had been primarily on the
60 sensor and data acquisition rather than data analysis or specific applications.  It could be argued that this
61 tendency has continued with Andersen et al. (2017), for example, stating that as of 2017 there was no
62 standardized accepted methodology for extracting surface points from green lidar point data alone –
63 although green and near-infrared lidar data are currently combined operationally to extract water surface

---

[1] A list of abbreviations is provided at the end of the article.

2

64    returns.  Nonetheless, recent review articles (Kashani et al. 2015; Kutser et al. 2020) suggest that

65    difficulties associated with extracting application-specific information from lidar point clouds are now

66    continually being addressed by the scientific community.  As lidar data processing/analysis research has

67    increased in recent years, new concepts for improving lidar sensors have also continued (e.g., Kinzel et al.

68    2021; Mandlburger et al. 2020; Mitchell and Thayer 2014).

69    This increased interest in lidar data analysis has been driven by a desire to decrease data acquisition costs

70    to better understand various phenomena.  For example, lidar data are proving to be particularly useful to

71    characterize benthic habitat.  Various analytical approaches have been explored: characterizing lidar

72    waveforms (Collin et al. 2008; Eren et al. 2018), using machine learning approaches (e.g., Pittman et al.

73    2009; Su et al. 2019), and classifying benthic habitat using variables that describe characteristics of lidar

74    pulses (Tulldahl and Wikstrom 2012) – 'soundings' in marine parlance.

75    The focus of this article is the use of airborne lidar data for shallow water bathymetry charting – defined

76    herein as water depths less than 20 m (although Jawak et al. 2015 noted that lidar can penetrate up to

77    60m under ideal conditions).  Much bathymetric depth work has focused on analysing the full lidar

78    waveform for a single spectral wavelength (see, for example, Pe'eri and Philpot 2007; Fernandez-Diaz et

79    al. 2014; Wang et al. 2015; Xing et al. 2019.) Waveform soundings have the potential to identify the water

80    surface and bottom due to increased reflectance from both.  Single wavelength waveform data are

81    operationally advantageous because they potentially decrease sensor complexity but are

82    disadvantageous because of increased data volumes compared to multi-wavelength systems that collect

83    point data.  Approaches to using waveform soundings for bathymetric mapping are varied and examples

84    include near-surface water modelling (Zhao et al. 2017), analysis of water column backscatter (Kinzel et

85    al. 2012; Nagle and Wright 2016), and a 'surface-volume-bottom' approach that provides a time-saving

86    closed-form solution (Schwarz et al. 2019).  The analysis of lidar point– rather than lidar waveform – data

87    has also received considerable attention (see, for example, Brzank et al. (2008), Yang et al. 2020) including

88    its combination with data from passive sensors (e.g., Dietrich 2017, Agrifiotis et al. 2019a).

89    Numerous researchers have examined ways of extracting bathymetry from such waveform data

90    algorithmically (e.g., Lyzenga et al. 2006, Pacheco et al. 2015, Li et al. 2019).  In such work, lidar is

91    sometimes used primarily as the reference data against which analytical methods are evaluated (e.g.,

92    Agrifiotis et al. 2019b).  In recent years, many such studies have examined various machine learning

93    techniques: neural networks (Liu et al. 2015), support vector machines (Misra et al. 2018; Wang et al.

94   2018), principal components analysis (Gholamalifard et al. 2013), partial least squares (Niroumand-Jadidi

95   et al. 2018), and random forests (Kogut and Weistock 2019).

96   The present research is focused on extracting shallow water bathymetry from lidar point clouds – i.e.,

97   identifying which lidar soundings represent the ocean bottom.  The approach adopted combines a density-

98   based algorithm developed for multi-beam echo sounder (MBES) sonar data with machine learning (ML)

99   techniques.  This methodological fusion is explored to overcome two considerable challenges.  First, lidar

100   data are collected from airborne platforms, resulting in a substantial number of soundings that represent

101   the ocean surface and near-surface.  Second, no ground-truth data are available for training ML models.

102   The latter difficulty was also recognized and addressed by Kerr and Purkis (2018), who developed a

103   workflow for optical data.  The former is suggestive of a weak bathymetric signal within a cloud of lidar

104   soundings.  The latter is particularly vexing because it creates a processing circularity: to determine which

105   lidar soundings represent ocean depth one needs at the least an initial depth estimate or, more ideally, a

106   *Bathy/NotBathy* designation for each sounding.  This article describes a method that overcomes both of

107   these difficulties and documents the results relative to a reference classification that is produced by

108   operationally adopted procedures.

109   **2.   Study Area and Lidar Data**

110   The airborne lidar data used for this work were captured by the United States National Oceanic and

111   Atmospheric Administration (NOAA) between April 22 and 25, 2016, in the vicinity of Key West, Florida

112   (24°33' N, 81°46' W).  Data were acquired by collecting lidar soundings over multiple overlapping flight

113   lines generally having a north-south orientation using a Riegel™ VQ-880-G sensor that employs a counter-

114   clockwise circular scan and a $20^0$ scan angle.  The nominal flying altitude of 400 m above mean sea level

115   results in an individual swath width of approximately 300 m and the pulse frequency of 45,000 pulses per

116   second provides a spatial density of approximately 10 soundings sq m$^{-1}$ for a single flight line.  The lidar

117   data were post-processed by NOAA by "cutting" the data from all flight lines into 500m-by-500m data

118   "tiles" aligned north-south and east-west with the Universal Transverse Mercator (UTM) projection.

119   Data for four tiles (Figure 1) were provided by NOAA in the format of the LAS data standard Version 1.4-

120   R13 (Point Data Record Format 6) (ASPRS 2013).   These tiles were selected because they are

121   representative of the range of sounding densities, depths, and ocean floor characteristics encountered in

122   operational shallow water bathymetric mapping (Table 1).  For convenience, the first five digits of each

123   tile's northing are employed as its identifier as well as a depth indicator – Shallow, Deep, Deeper, or

124   Deepest.  The overlap in flight lines produces a combined average sounding density between 13 and 30

125 returns m$^{-2}$, although sounding density varies across each tile and is considerably higher where flight lines

126 overlap.

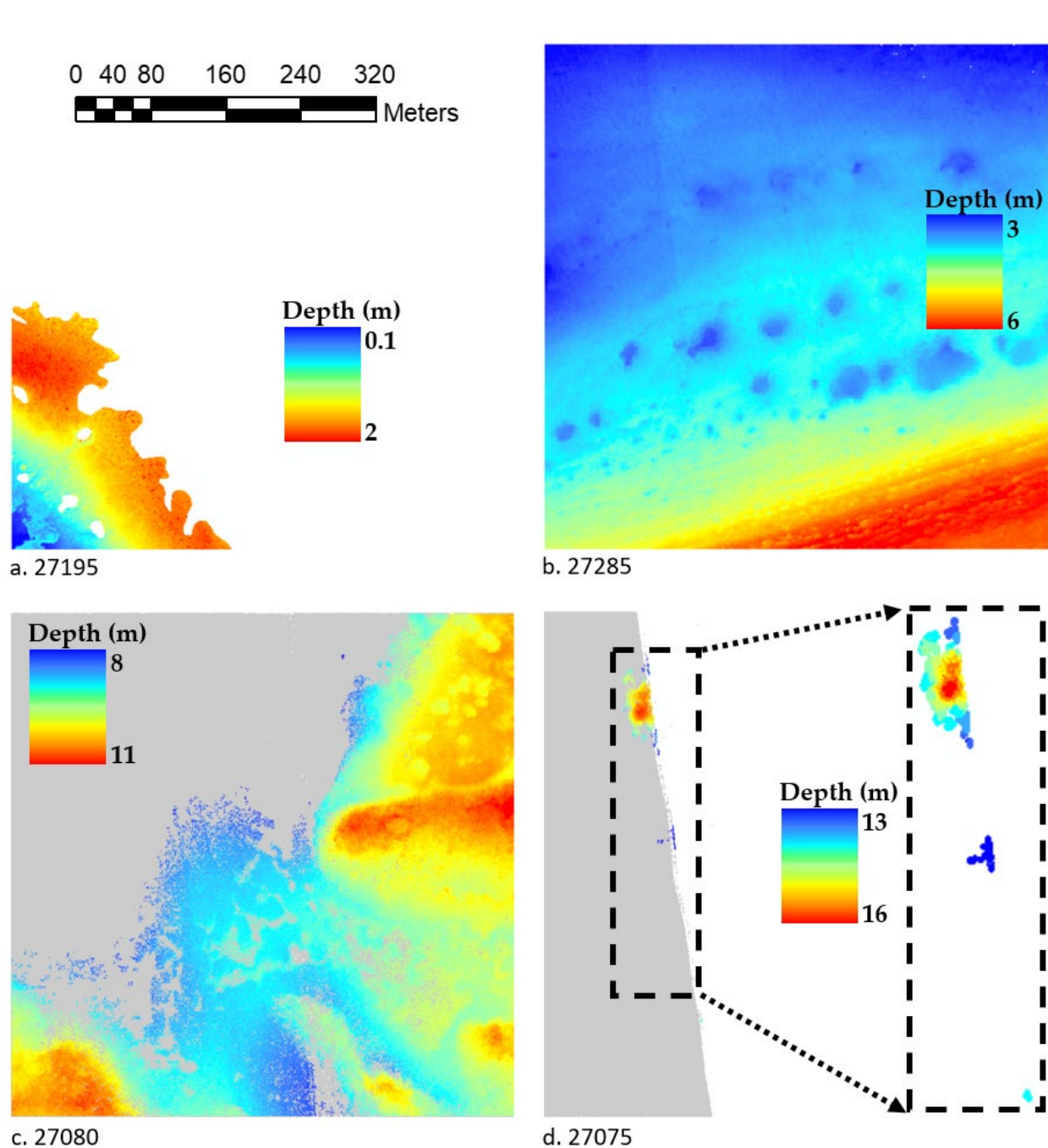127



a. 27195

b. 27285

c. 27080

d. 27075

128
129 Figure 1. Depth maps (1m pixels) for the four tiles based on depth determined by NOAA. White areas
130 have no usable data. Gray areas have usable data, but no soundings were identified as bathymetry by
131 NOAA. Due to sparseness of NOAA-identified bathymetry on the deepest tile (27075; Fig. 1d), an
132 enlargement of the area containing bathymetry is shown with bathymetric soundings accentuated and
133 gray background removed.

134 Table 1. Descriptive information about the data tiles employed in this study.

| Identifier (Northing) | Relative Depth | Description | Area (km²) | Approx. MSL depth range (m) | Total Soundings (million) | Mean return density (pts/m²) | % Bathymetry | Number of flight lines |
|---|---|---|---|---|---|---|---|---|
| 27195 | Shallow | Shallow area including some mangrove swamps | 3 | 0 to 2 | 0.6 | 27.6 | 78 | 5 |
| 27285 | Deep | Gradual slope with a few scattered mounds about 1 m tall | 25 | 3 to 6 | 7.6 | 30.4 | 76 | 7 |
| 27080 | Deeper | Gradual slope cut by relatively shallow channels; the northwest is poorly classified | 25 | 8 to 11 | 3.7 | 14.8 | 21 | 7 |
| 27075 | Deepest | Depth mostly beyond limit of lidar penetration except for mound in northeast and isolated points on eastern edge. | 7.5 | 13 to 16 | 0.9 | 13.3 | 0.4 | 2 |

135

136 Attached to each sounding are its geographic (UTM) coordinates, depth, time of acquisition, and a variety
137 of metadata that we term "sounding attribute data" (SAD; Table 2). Lidar depth is expressed in meters
138 relative to mean sea level determined using NOAA's VDATUM tool (https://vdatum.noaa.gov/). Sounding-
139 based SAD are either acquired by the lidar instrument or were derived post-acquisition. Also provided
140 were Smoothed Best Estimate of Trajectory (SBET) data. These are produced by the Applanix software by
141 post-processing pulse return data from each flight line using a proprietary method based on a tightly
142 coupled extended Kalman filter. SBET data have had noise removed to describe the most likely airplane
143 position and orientation at 200 Hz. Flight path and orientation consistency are described in the SBET data
144 by the standard deviations for the x, y, and z location of the plane and its yaw, pitch, and roll extracted
145 from the Kalman filter's post-observation covariance matrix. SBET values were assigned to individual
146 soundings by matching time of acquisition.

147 An additional variable was created to characterize platform stability at the moment of data acquisition. It
148 was observed that the crenularity – i.e., the deviation from a straight line -- of the margin of soundings of
149 individual flight lines varied along the flight line (Figs. 2a and 2b). We hypothesized that this crenularity
150 reflected local wind conditions that may in turn impact surface water conditions and lidar reflectance
151 characteristics. To quantify this, sounding cloud 'edge points' were identified algorithmically along the
152 length of the flight path (Fig. 2a). The two end soundings were considered 'corner' soundings and the
153 equation of the straight line between them calculated (Fig. 2b). The orthogonal distance from each edge

154    sounding to the straight line was determined (Fig. 2c) and the absolute value of this deviation was assigned

155    to each sounding by matching time of acquisition.  This variable is termed *abs_devia*.

156    Table 2. Depth and sounding attribute data (SAD) employed in this study for machine learning (ML)
157    modelling.  Variable names are *italicised* throughout the article.

| SAD Type | Nature | Variable (*Name*) |
|---|---|---|
| **Depth** | ***Depth*** | Depth as provided via lidar post-processing (*depth* in m) |
| **Sounding-based** | ***Pulse-specific*** | • Intensity of sounding return (*intensity*; 16-bits – i.e., maximum value is 65536) <br> • Number of soundings (*numreturns*) <br> • Sounding number from a given lidar pulse (*return_no*) <br> • First sounding of many (*first_of_many*; 0 or 1) <br> • Last sounding of many (*last_of_many*; 0 or 1) <br> • Last sounding (*last*; 0 or 1) <br> • Scan direction (*scan_direct;* -1 (backwards) or +1 (forward)) <br> • Azimuth from airplane to pulse (*azim2plse*; $0^0$ to $360^0$ in decimal degrees) <br> • Incident scan angle corrected for yaw, pitch, and roll (*inciangle*; recorded in decimal degrees) <br> • Difference between pulse direction and airplane heading (*plse_frm_hdng*; $0^0$ to $90^0$ decimal degrees) |
| **Airplane stability** | ***SBET*** | • Aircraft positional – sum of standard deviations of x, y, and z (*stdXYZ*) <br> • Aircraft platform – sum of standard deviations of Yaw, Pitch, and Roll (*stdYwPtRl*) <br> • Deviation from flight path – see Figure 2 and text for explanation (*abs_devia*) |

158

159    Most of the SAD variables can be considered 'direct features' (Höfle and Rutzinger 2011) that are

160    measured, although *abs_*devia would be considered an 'indirect feature' as it is derived post-acquisition.

161    The SAD variables in Table 2 were retained for analysis because of their documented or hypothesized

162    impact on light reflectance directly and bathymetric signal indirectly.  Examples of their documented

163    impact include *intensity* (Schmidt et al. 2012), *abs_devia* and SBET variables as surrogates for surface

164    waves (Westfeld et al. 2017; Maas et al. 2019; Lowell et al. 2021), and *inciangle* (Birkeback et al. 2018;
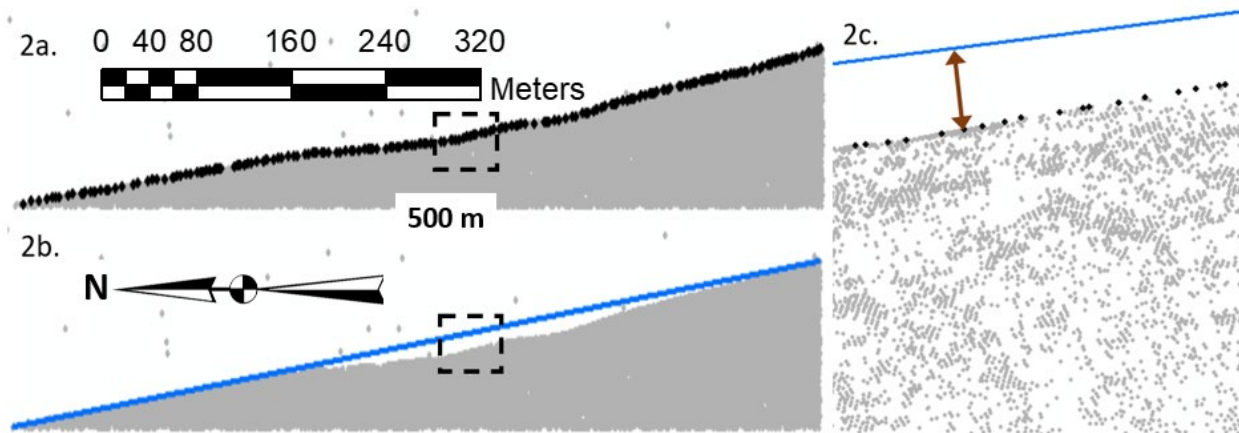
165    Okhrimenko and Hopkinson 2020).

166

Figure 2. Extraction of orthogonal deviations from a single flight line (the western edge of the deepest tile 27075; Fig. 1d). To conserve space, North points to the left. a) Lidar point cloud (gray) and derived edgepoints (black). The dashed box is the area of enlargement in Fig. 2c. b) "Corner-to-corner" straight flight path (blue line). c) Orthogonal deviation of a single edge point from the corner-to-corner flight path edge of lidar point cloud.

Finally, also available for each sounding is the *Bathy/NotBathy* classification produced by NOAA using in-house methods. The LAS data standard classes of interest herein are 'Bth', 'Unc', and 'LP-Nz' that generally represent bathymetry, water surface, and water column noise, respectively. For the current work, these were condensed into two classes – *Bathy* ('Bth' only) and *NotBathy* ('Unc' and 'LP-Nz'). This *Bathy/NotBathy* classification is used only as the reference classification against which the results of the method developed are compared – i.e., it is not used in the method developed. Moreover, although the NOAA *Bathy/NotBathy* classification is the most authoritative available and is an appropriate standard for comparison since it is used operationally, it is not 'ground truth' produced via direct measurement or observation.

For analysis, notable in the data tiles employed is the spatial distribution of *Bathy* – i.e., the not-gray points in Figure 1. Tile 27195/Shallow (Fig. 1a) is the shallowest tile, is located in an area of mangrove swamps, and all of it except the southwest area has been classified by NOAA as being above sea level. For such areas, NOAA creates a data exclusion mask so that only soundings from aquatic areas are classified as *Bathy/NotBathy*. For consistency, the same practice is adopted herein and only data from the colored area shown in Fig. 1a are employed in subsequent analyses. Also notable – and representative of real-world conditions – are the incomplete *Bathy* coverages of Tiles 27080/Deeper (Fig. 1c) and 27075/Deepest (Fig. 1d). This results from increasing depths that ultimately exceed the depth limit of lidar penetration. It is most pronounced for Tile 27075/Deepest on which only 0.4% of soundings are *Bathy* (Table 1).

192 **3.   Procedures**

193   Figure 3 is a schematic showing the procedural flow of the work undertaken.  This is explained below and

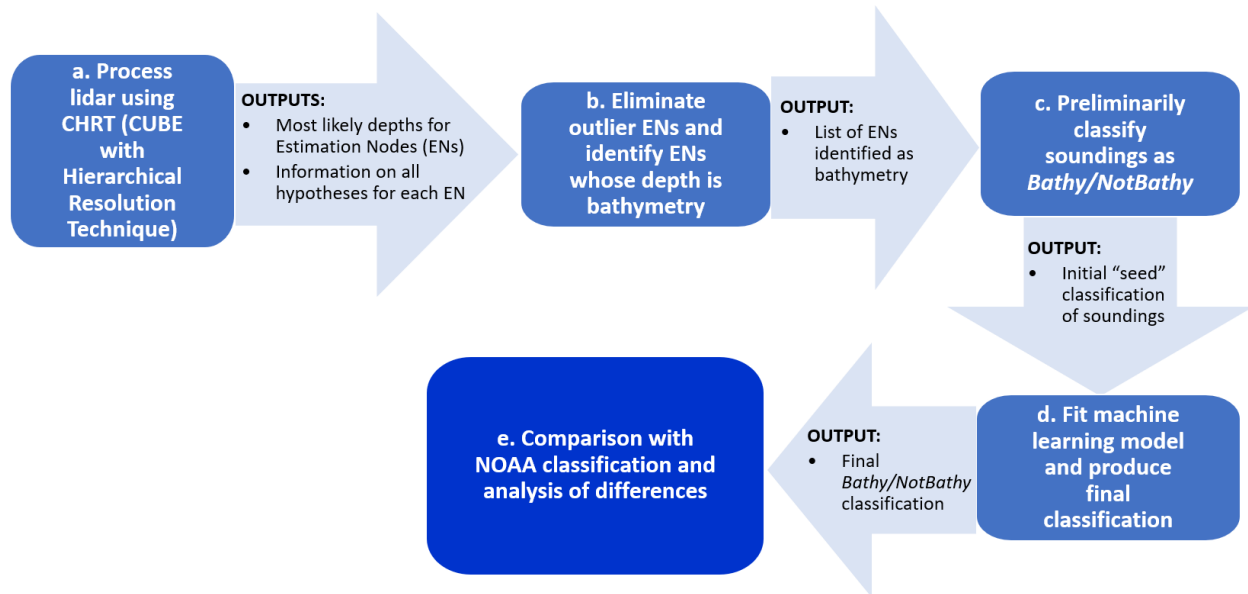194   was applied to each data tile individually.

195



196
197   Figure 3. Schematic of data processing flow.  Letters refer to parts of section 3 of the paper; the final step

198   ("e.") is not further described in section 3.

199   **3a.        Process lidar using density-based/CHRT algorithm**

200   This process identifies the most likely depth (MLD) for a north-south/east-west grid of 'estimation nodes'

201   (ENs) established over each lidar data tile.  This information is produced by processing the lidar point cloud

202   data through a density-based algorithm as described below.

203   The algorithm employed in this study is CHRT (CUBE with Hierarchical Resolution Technique; Calder and

204   Rice 2017) which is a modification of CUBE (Combined Uncertainty and Bathymetry Estimator; Calder and

205   Mayer 2003).  CHRT is incorporated into many software packages that are widely used operationally and

206   scientifically for processing MBES sonar data (Lecours et al. 2016).  Its scientific use has also been extended

207   into other applications such as benthic habitat mapping (Calvert et al. 2015).

208   CHRT establishes a grid of ENs across an area of interest with the spacing of the grid determined by the

209   density of the soundings.  Given the non-rectangular nature of the spatial coverage for some tiles in this

210   study (see Figure 1), some ENs on a grid have no associated lidar soundings and are removed from further

9

211    analysis.  Similarly, as will be explained subsequently, others have aberrant MLD values due to data

212    anomalies and sparseness; these are also removed using outlier analysis.

213    To estimate the MLD for a single EN, the 'neighboring' soundings for the EN are identified.  An EN's

214    'neighboring' soundings are those within a geographic radius defined by the grid spacing for a tile.  The

215    radius used is the Euclidean distance to the 'pixel corner' defined by a point equidistant from four adjacent

216    ENs (Figure 4).  Thus 50-60% of soundings are neighbors of two ENs depending on local sounding density.
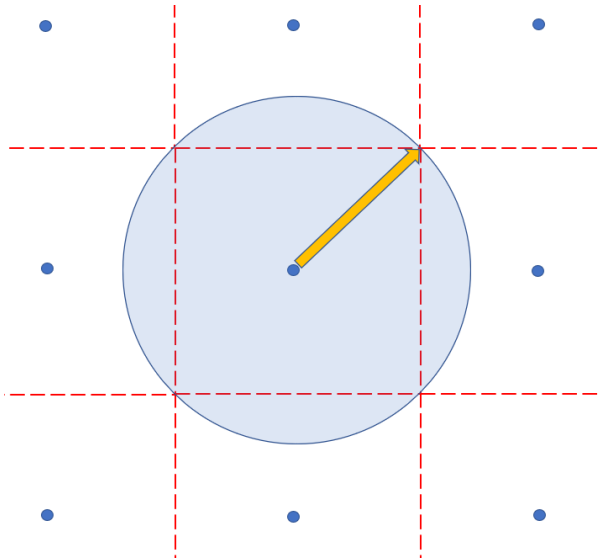


217
218    Figure 4. Search radius to determine "neighboring" pulse returns.  Blue points are estimation nodes (ENs).
219    The orange arrow defines the neighbor radius.

220    For each EN, its neighboring soundings are progressively ingested with the first sounding defining an initial

221    depth 'hypothesis.'  A variety of hyperparameters or 'tuning' parameters establish initial thresholds for

222    determining if two soundings represent different depths; default parameters are based on user

223    experience with the location and variability of depth frequency distributions for various depth conditions.

224    The depth of the second sounding ingested is evaluated against the hyperparameters to determine if it

225    also 'belongs to' the first depth hypothesis, or if its depth is 'different enough' to be considered a new

226    hypothesis.  This process continues until all neighboring soundings have been ingested and one or more

227    depth hypotheses have been developed and characterized.  As the process progresses, frequency

228    distributions for each hypothesis are produced and their characteristics – rather than the initial

229    hyperparameters -- increasingly control the assignment of newly ingested soundings to existing

230    hypotheses, or to the creation of a new hypothesis.  After all soundings have been ingested,

231    disambiguation rules determine which hypothesis represents the MLD for the EN.  A simplistic example is

232    that the hypothesis with the deepest mean depth is considered the MLD although such a rule ignores

233 factors such as turbidity in the water column or the number of soundings in the deepest hypothesis;

234 current disambiguation rules generally identify the hypothesis having the greatest number of soundings

235 as most likely.

236 Table 3 provides summary information about the ENs for each tile. One notable point is that the variability

237 (standard deviation) of MLDs for tile 27080/Deeper is considerably larger than for other tiles. This reflects

238 both a greater variability in geomorphometry for this tile and the large area in the northwest beyond the

239 range of lidar penetration (Fig. 1c). Note that because each tile is processed individually (here and

240 operationally) and the density of the EN grid on each tile depends on its sounding density, EN density

241 varies across tiles. This causes edge artifacts when combining adjacent tiles into a seamless map;

242 procedures for doing this are beyond the scope of this study.

243 Table 3. Estimation node (EN) information after removal of outliers. (See text for explanation.)

| Identifier (Northing) | Relative Depth | Grid Spacing (m) | EN Grid: Rows*Cols | ENs used for analysis[1] | Mean hypotheses per EN | Mean soundings per hypothesis | MLD[2] range (m) | Mean MLD (m) | MLD standard deviation (m) |
|---|---|---|---|---|---|---|---|---|---|
| 27195 | Shallow | 12.4 | 20*18 | 188 | 4.8 | 3827 | 0 to 1 | 0.7 | 0.2 |
| 27285 | Deep | 1.6 | 308*308 | 94853 | 4.6 | 125 | 2 to 7 | 4.5 | 0.7 |
| 27080 | Deeper | 3.0 | 167*167 | 27823 | 6.3 | 32 | 1 to 10 | 3.5 | 2.3 |
| 27075 | Deepest | 1.9 | 267*120 | 18446 | 5.2 | 82 | 1 to 18 | 1.5 | 0.5 |

244 [1]Estimation Nodes after removal of no-data and outlier estimation nodes.

245 [2]Most likely depth (MLD).

246

247 **3b      Eliminate outliers and identify ENs whose MLD is bathymetry**

248 A two-phase outlier screening process is employed; Table 4 provides information on the results of this

249 screening.

250 Table 4. Information about estimation node (EN) outlier screening.

| Identifier (Northing) | Relative Depth | Total Grid ENs | ENs w/o Soundings | MD[1] Outliers | Beyond Lidar Penetration MLDs[2] | ENs Analysed |
|---|---|---|---|---|---|---|
| 27195 | Shallow | 360 | 136 | 2 | 37 | 185 |
| 27285 | Deep | 94864 | 11 | 1232 | 0 | 93621 |
| 27080 | Deeper | 27889 | 0 | 318 | 14 | 27557 |
| 27075 | Deepest | 32040 | 13514 | 288 | 34 | 18204 |

251 [1]Mahalanobis Distance.

252 [2]Most Likely Depth (MLD).

253 First, 12 variables associated with each EN's hypotheses are used to calculate Mahalanobis distances

254 (MDs; Mahalanobis 1936) for each EN.  Examples of such variables are the number of hypotheses, total

255 sounding, the number of soundings associated with the MLD hypothesis and non-MLD hypotheses, and

256 the standard deviation of the depth of soundings associated with the MLD and non-MLD hypotheses.

257 Prior to calculating the MDs, variables are normalized between 0 and 100 using max-min normalization.

258 ENs are eliminated from subsequent analysis if their MD is in the outer 0.1% of the frequency distribution

259 – i.e., their MD is more than approximately 3.3 standard deviations from the mean MD.

260 Second, airborne lidar cannot penetrate below certain ocean depths.  Examination of depth frequency

261 distributions across all tiles suggested that for the area studied, lidar could not penetrate below a depth

262 of 20 m.  Hence ENs whose MLD depth was greater than 20 m and that had not already been removed by

263 the MD outlier analysis are eliminated as ''Beyond Lidar Penetration' MLDs.

264 MLD frequency distributions for the four tiles were highly irregular – e.g., not clearly normal or bi-modal

265 – generally reflecting a separation of ocean surface and ocean bottom.  Hence *k*-means clustering

266 (Steinhaus 1957, McQueen 1967) is applied to the MLD of the ENs retained for analysis to separate them

267 into two classes.  Because a single variable – MLD -- is used in this clustering, this is equivalent to

268 separating a frequency distribution along a single axis. The cluster having the greatest difference between

269 its MLD and the average depth of all other hypotheses – i.e., the 'non-MLD' hypotheses -- is assumed to

270 contain some EN hypotheses that are 'definitely' bathymetry.  The other cluster has a smaller difference

271 between the mean MLD and the mean depth of non-MLD hypotheses suggesting that both represent the

272 ocean surface.  Note that not all ENs will represent ocean floor or surface, but some will represent the

273 water column.  This is most likely to be problematic where water column soundings are more prevalent

274 than ocean floor soundings – i.e., in highly turbid waters or beyond the limits of lidar depth penetration.

275 The mean and standard deviation for the MLDs of the cluster identified as 'definitely' containing

276 bathymetry are used to define the bathymetry MLD confidence interval for the MLDs for all ENs.  The

277 shallower MLD limit of the interval is the one-sided 99.9% confidence limit whereas the deeper MLD limit

278 is the one-sided 95% confidence limit.  These 'imbalanced limits' were found to address the irregularly

279 shaped bathymetry frequency distributions across all tiles better than equal 'shallower/deeper' MLD

280 confidence limits.

281 The MLD bathymetry confidence interval is used to classify the MLDs of all ENs as *Bathy* or *NotBathy*.  That

282 is, the MLD hypotheses of all ENs contained within the bathymetry confidence interval are classified as

283   *Bathy*.  All other hypotheses, even including those whose mean depth falls in the bathymetry depth
284   interval but are not the MLD for their EN, are classified as *NotBathy*.

285   Alternatives to this classification rule were evaluated, including clustering on all hypotheses rather than
286   on MLD hypotheses only, classifying as *Bathy* all hypotheses – MLD and non-MLD -- whose mean depth
287   fell in the bathymetry depth interval, and using the range instead of the standard deviation to define the
288   bathymetry confidence interval.  None performed as well as the clustering approach and classification
289   rule adopted.

290        **3c        Preliminarily classify pulse returns as *Bathy* or *NotBathy***
291   The neighboring soundings for each EN are classified as *Bathy* if they are associated with the MLD and the
292   MLD has been classified as *Bathy*; otherwise they are classified as *NotBathy* (Table 5).  Soundings that are
293   neighbors of two ENs whose *Bathy* or *NotBathy* classification agrees are assigned to the agreed class.
294   Two-neighbor soundings whose classifications do not agree are termed 'mixed' and are assigned to *Bathy*.
295   Tile 27195/Shallow with the least separation between ocean surface and ocean depth had the largest
296   percentage of mixed soundings and Tile 27075/Deepest had the smallest (Table 5).  This assignment
297   scheme for mixed soundings has the effect of decreasing the number of false negatives (FNs) – undetected
298   bathymetry – which is a more serious error than a false positive (FP) – erroneously labelled bathymetry –
299   in nautical chart production.  Assigning mixed soundings to *Bathy* also was found to improve the skill of
300   the subsequently fitted machine learning model.

301   Table 5. *Bathy/NotBathy* classification information for soundings that were neighbors of two estimation
302   nodes (ENs).

| Identifier (Northing) | Relative Depth | Soundings Neighboring Two ENs | Pure Bathy Soundings | Pure NotBathy Soundings | "Mixed" Soundings | % "Mixed" |
|---|---|---|---|---|---|---|
| **27195** | Shallow | 307500[1] | 138500 | 85700 | 83300 | 27 |
| **27285** | Deep | 4246400 | 2307400 | 1547700 | 391300 | 9 |
| **27080** | Deeper | 2081400 | 600900 | 1343100 | 137400 | 7 |
| **27075** | Deepest | 542100 | 3100 | 538400 | 600 | <1 |

303   [1]All values rounded to nearest 100.

304        **3d        Fit a machine learning model to produce a final *Bathy/NotBathy* classification for all**
305   **soundings**
306   At this point a preliminary or 'seed' classification has assigned each sounding to the *Bathy* or *NotBathy*
307   class.  However, only soundings associated with the MLD hypothesis will have been classified as *Bathy*.
308   This is problematic in areas where all hypotheses – MLD and others – are representative of bathymetry.

13

309    This occurs on tiles where bathymetry is commonplace (27195/Shallow and 27285/Deep; Table 1), as well

310    as in areas where bathymetry is locally concentrated such as the northeastern area of 27075/Deepest

311    (Fig. 1d). Hence this assignment is a 'preliminary' or 'seed' classification that can be considered

312    conservative – i.e., soundings classified as *Bathy* have an extremely high 'true probability' of being *Bathy,*

313    but the classification likely contains a high number of FN errors.

314    Therefore, to produce a final classification, extreme gradient boosting (XGB) (Friedman 2001) is used to

315    fit a model using the seed classification. XGB is a decision-tree machine-learning technique that

316    progressively fits numerous simple or 'shallow' models/trees. Each successive tree is fitted with a focus

317    on the worst-predicted observations of the previous tree. Once statistical convergence or the maximum

318    number of trees is achieved, a composite XGB model is produced using a 'majority vote' approach. Lowell

319    et al. 2021) have demonstrated that the SAD employed (Table 2) contain a substantial amount of

320    bathymetric signal. Specifically, machine learning models that used the variables in Table 2 as

321    independent variables and NOAA's *Bathy/NotBathy* as the dependent variable produced $R^2$ values

322    between 0.61 and 0.99 and global classification accuracies between 90% and 99.9%.

323    Hence, an XGB model that uses the seed CHRT/clustering-based classification as the dependent variable

324    and the variables in Table 5 as predictor variables is fitted for each tile. This model is then used to estimate

325    *p(Bathy)* – i.e., the probability of each pulse return being *Bathy* for each sounding. Soundings are classified

326    as *Bathy/NotBathy* by applying to the *p(Bathy)* values a probability decision threshold (PDT) that equalizes

327    the true positive (*Bathy*) rate (TPR) and true negative (*NotBathy*) rate (TNR) rather than the conventional

328    PDT of 0.50; we term this alternative PDT the "optimal decision threshold" or "ODT." Lowell et al. (2021)

329    demonstrated that the use of the ODT mitigates the impacts of the accuracy of a class – *Bathy* or *NotBathy*

330    – comprising a strong majority of soundings being maximized at the expense of the accuracy of the

331    minority class. Problems associated with applying a conventional PDT of 0.50 were notable for all tiles,

332    but especially for 27075/Deepest on which only 0.4% of pulse returns were identified by NOAA as being

333    *Bathy*.

334    **4. Analysis, Results, and Discussion**

335    Initially assessed are 1) how well the 'preliminary/seed' *Bathy/NotBathy* classification derived from

336    clustering the EN MLDs performs relative to the NOAA reference classification and 2) if the final XGB

337    model-based classification improves the preliminary seed classification relative to the NOAA reference
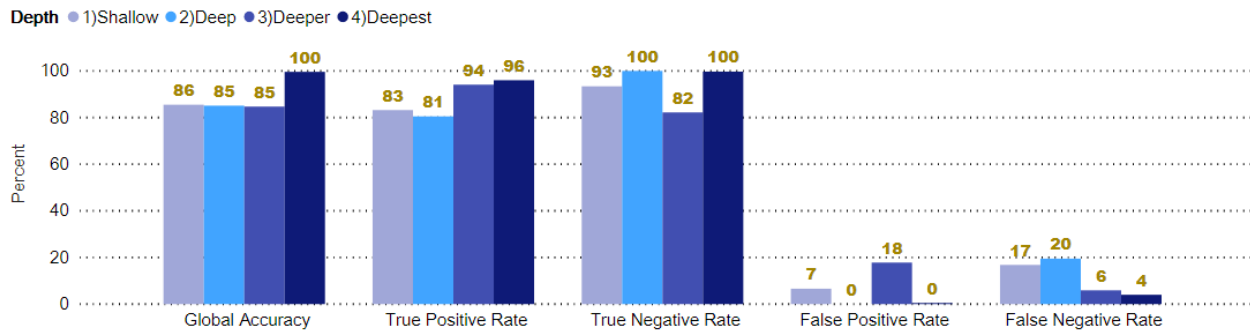
338   classification.  Methods that can be used for continuous improvement of classification methodology are

339   subsequently presented.

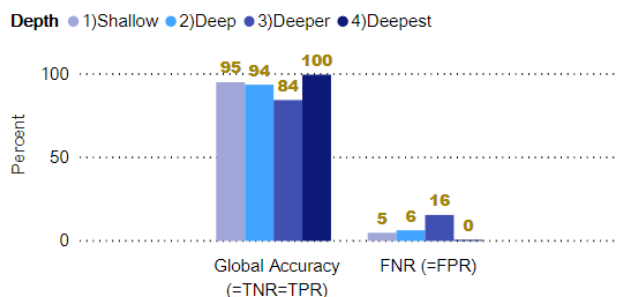### 4a. Classification accuracies

341   Recall that the preliminary seed classification produced by clustering EN MLDs from the CHRT algorithm

342   is considered a conservative classification because only the most certain soundings are classified as *Bathy*.

343   Nonetheless, because it will be subsequently refined using an XGB model, it only needs to be 'sufficiently

344   accurate' that the XGB model will be able to detect and describe underlying relationships between

345   bathymetry and the SAD variables (Table 2).  If this occurs, the XGB model should be able to identify

346   soundings not classified as *Bathy* initially, but that have a high *p(Bathy)* nonetheless.  Subsequently

347   reclassifying all soundings based on the *p(Bathy)* values should thus expand the number of soundings

348   correctly classified as *Bathy* or *NotBathy*.  A truly ideal outcome would be that the preliminary seed

349   classification is identical to the NOAA reference classification and thus does not require additional

350   processing.

351   Figure 5a suggests that the preliminary seed clustering classification relates strongly to the NOAA

352   reference classification.  Global accuracy – or more precisely 'agreement between the two' – is the

353   percentage of all soundings that are correctly classified as *Bathy* or *NotBathy*; it is at least 85% for all tiles.

354   Moreover, the percent of correctly classified *Bathy/NotBathy* soundings – true positives (TPs) and true

355   negatives (TNs), respectively – is 81% or more for all tiles.  Similarly, the percentage of false positives (FPs;

356   *NotBathy* soundings incorrectly labelled *Bathy*) and false negatives (FNs; *Bathy* soundings incorrectly

357   classified as *NotBathy*) is 20% or lower for all tiles.

358

a. NOAA Reference vs Preliminary Seed Classification

b. NOAA Reference vs Final (XGB) Classification

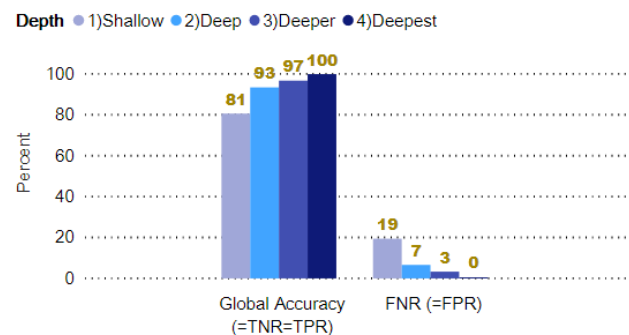c. Final (XGB) vs Prelim. Seed Classification

Figure 5. Classification accuracy/comparisons for various classification schemes.

The ability of an XGB model to improve on the preliminary classification – i.e., to harmonize it with the NOAA classification -- can be assessed by comparing Figs. 5a and 5b. Readers are reminded that to classify *p(Bathy)* as *Bathy/NotBathy*, the ODT that equalizes the TNR and TPR was employed throughout. Since the FNR is 100 minus the TPR, the FPR is 100 minus the TNR, and the TPR and TNR are equal, the FNR and FPR are equal. The use of the ODT also has the impact of making the global accuracy equal to the TNR and TPR – e.g., if the classification of *NotBathy* soundings is 80% correct and the classification of *Bathy* soundings is 80% correct, global accuracy for all soundings must also be 80%.

Figure 5b suggests that the XGB model improved the initial cluster-based 'seed' classification. Through the use of an XGB model, global accuracy improved for all tiles except 27080/Deeper for which it decreased by a single percent (85% to 84%). The TPR follows a similar pattern, although the TNR did decrease for Tile 27285/Deep (from 100% to 94%). Of greatest real-world interest is that the FNR has dropped considerably for all tiles except 27080/Deeper. In operational practice, a decrease in FNs not only means improved navigational safety but also considerable cost savings. That is, because FNs are the most serious error for nautical navigation, considerable human-time is spent verifying FNs. Hence reducing the number of FNs decreases the time spent on manual editing. Thus except for the 27080/Deeper tile, the ML model has improved classification in a way that benefits operational workflows

16

377 and improves navigational safety.  We note a considerable portion of the 27080/Deeper tile (northwest

378 of Figure 1c) exceeds the depth of lidar penetration which undoubtedly impacts the accuracy of identifying

379 *Bathy* soundings.

380 Results can potentially be further improved by better understanding the XGB model fitted to the binary

381 *Bathy/NotBathy* cluster-based "seed" classification using the variables in Table 2 as the independent

382 variables.  The classification accuracy of the model on this preliminary classification– i.e., the classification

383 used to fit the model rather than NOAA's reference classification -- is at least 81% for all tiles (Fig. 5c).  The

384 goodness-of-fit/explanatory power of the model can also be evaluated by calculating an $R^2$ for binary

385 dependent variables (McFadden1974) that is conceptually equivalent to, and interpreted in the same

386 manner as, the more familiar $R^2$ value associated with linear regression.  The $R^2$ values for individual

387 models are relatively high (Table 6) – particularly considering the large number of soundings (at least

388 600,000; Table 1) used to fit the models.  The number of variables with an 'importance value' (a measures

389 of a variable's contribution to an XGB model) greater than zero (0) was at least 11 for all tiles and the five

390 most important variables contained at least 95% of the total importance for all tiles except 27195/Shallow.

391 Coupled with the fact that other than *depth*, an inconsistent variety of SAD variables were important, the

392 information that provides discrimination between *Bathy* and *NotBathy* soundings appears to be

393 distributed among a suite of variables specific to each tile;  XGB as a model development technique is able

394 to accommodate this variability.

395 These findings are potentially most relevant for the 27195/Shallow tile.  Its relatively low cumulative

396 importance of the five most important variables (0.85) suggests that for shallow areas where the distance

397 between the ocean surface and ocean floor is less than the noise in the lidar sounding cloud, *depth*

398 contains a smaller proportion of the information the provides discrimination between *Bathy* and *NotBathy*

399 soundings than it does for deeper areas.  Finally, SBET variables were among the five most important

400 variables in two of the models and *abs_devia* was present in one suggesting that both SAD associated with

401 individual soundings and SAD that describe flight path and airplane stability contain information that

402 provides discrimination between *Bathy* and *NotBathy* soundings.

403 Table 6. Information about XGB models.

| Identifier (Northing) | Relative Depth | R-squared[1] | Number of Important Variables | Five most important variables[2] | Cumulative importance of the five most important variables |
|---|---|---|---|---|---|
| **27195** | Shallow | 0.48 | 12 | *depth, last, first_of_many, stdYwPtRl, inciangle* | 0.85 |

| 27285 | Deep | 0.79 | 14 | *depth, return_no, last_of_many, intensity, plse_frm_hdng* | 0.95 |
|--------|------|------|----|---|------|
| 27080 | Deeper | 0.87 | 12 | *depth, num_returns, return_no, last_of_many plse_frm_hdng* | 0.99 |
| 27075 | Deepest | 0.97 | 11 | *depth, plse_frm_hdng, abs_devia, scan_direct, stdXYZ* | 0.99 |

404    [1]McFadden's (McFadden 1974) pseudo $R^2$ which cannot be tested for statistical significance.

405    [2]In descending order of importance.

406    **4b. Continuous Improvement.**

407    Because neither the NOAA classification nor the XGB classification developed can be considered ground

408    'truth' that results from direct measurement, being able to characterize differences between the two is

409    useful for improving the NOAA classification, the final XGB classification, or both.  Two methods were

410    developed to provide such information.

411    The first method focusses on 'feature' or 'statistical' space and entails comparing *p(Bathy)* values from

412    the XGB model with the NOAA *Bathy/NotBathy* classification using logistic regression.  The approach is

413    comparable to binning the soundings by *p(Bathy)* values and then determining if the proportion of

414    soundings classified as *Bathy* by NOAA in each bin equals the bin mid-point class value.  The logistic

415    regression approach employed, however, provides information along the continuum of *p(Bathy)* values

416    without requiring an arbitrary number of bins.  In this approach, NOAA's *Bathy/NotBathy* classification is

417    used as the dependent variable and the following logistic equation is fitted:

418   
$$p' = \left(1 + e^{\left(-(b_0 + b_1 L)\right)}\right)^{-1} \tag{1}$$

419    where L is:

420   
$$L = ln\left(\frac{p}{(1-p)}\right) \tag{2}$$

421    and *p* is the *p(Bathy)* estimated by the XGB model.  For each tile, if the NOAA classification and *p(Bathy)*

422    values from the XGB model are identical over the entire *p(Bathy)* range of 0.0 to 1.0, $b_0$ and $b_1$ in Equation

423    (1) will be 0.0 and 1.0, respectively.  Furthermore, $R^2$ for Equation 1 will be 1.0 with an associated log-

424    likelihood *p* that is infinitesimally small.  Such a 'logistic agreement model' was fitted for each tile (Table

425    7).

426    Table 7. Information on logistic agreement models.

| Identifier (Northing) | Relative Depth | R-squared[1] | log-likelihood $p$ | $b_0$[2] | $b_1$[3] | n |
|---|---|---|---|---|---|---|
| 27195 | Shallow | 0.85 | <0.001 | 1.1* | 2.05* | 576,000 |
| 27285 | Deep | 0.79 | <0.001 | 19.3* | 2.83* | 7,599,000 |
| 27080 | Deeper | 0.52 | <0.001 | -1.3* | 0.60* | 3,706,000 |
| 27075 | Deepest | 0.77 | <0.001 | -2.3* | 0.70* | 983,000 |

427   [1]McFadden's pseudo $R^2$ that cannot be tested for statistical significance.

428   [2] * signifies the intercept value is significantly different from 0.0 at $a$=0.001.

429   [3] * signifies the slope value is significantly different from 1.0 at $a$=0.001.

430   The relatively high $R^2$ values and low log-likelihood $p$ values for the logistic models (Table 7) suggest a

431   strong and significant relationship between the *p(Bathy)* produced by the XGB model fitted on the

432   preliminary classification and the NOAA *Bathy*/*NotBathy* classification.   However, for all models the

433   intercepts ($b_0$) and slopes ($b_1$) are significantly different from 0.0 and 1.0, respectively.

434   To assess (dis)agreement over the entire probability range, the logistic agreement models of Table 7 can

435   be displayed graphically by plotting *p'* vs. *p* over the interval {0,1}.  This was done using the ODT that is

436   specific to each tile – i.e., by stretching *p(Bathy)* values in the range {0, ODT} to the interval {0.0, 0.5}, and

437   stretching *p(Bathy)* values in the range {ODT, 1.0) .0 to the interval {0.5, 1.0}.  Note that this segmenting

438   of the probability range is the cause of the graphical discontinuities present at a *p(Bathy)* value of 0.5 for
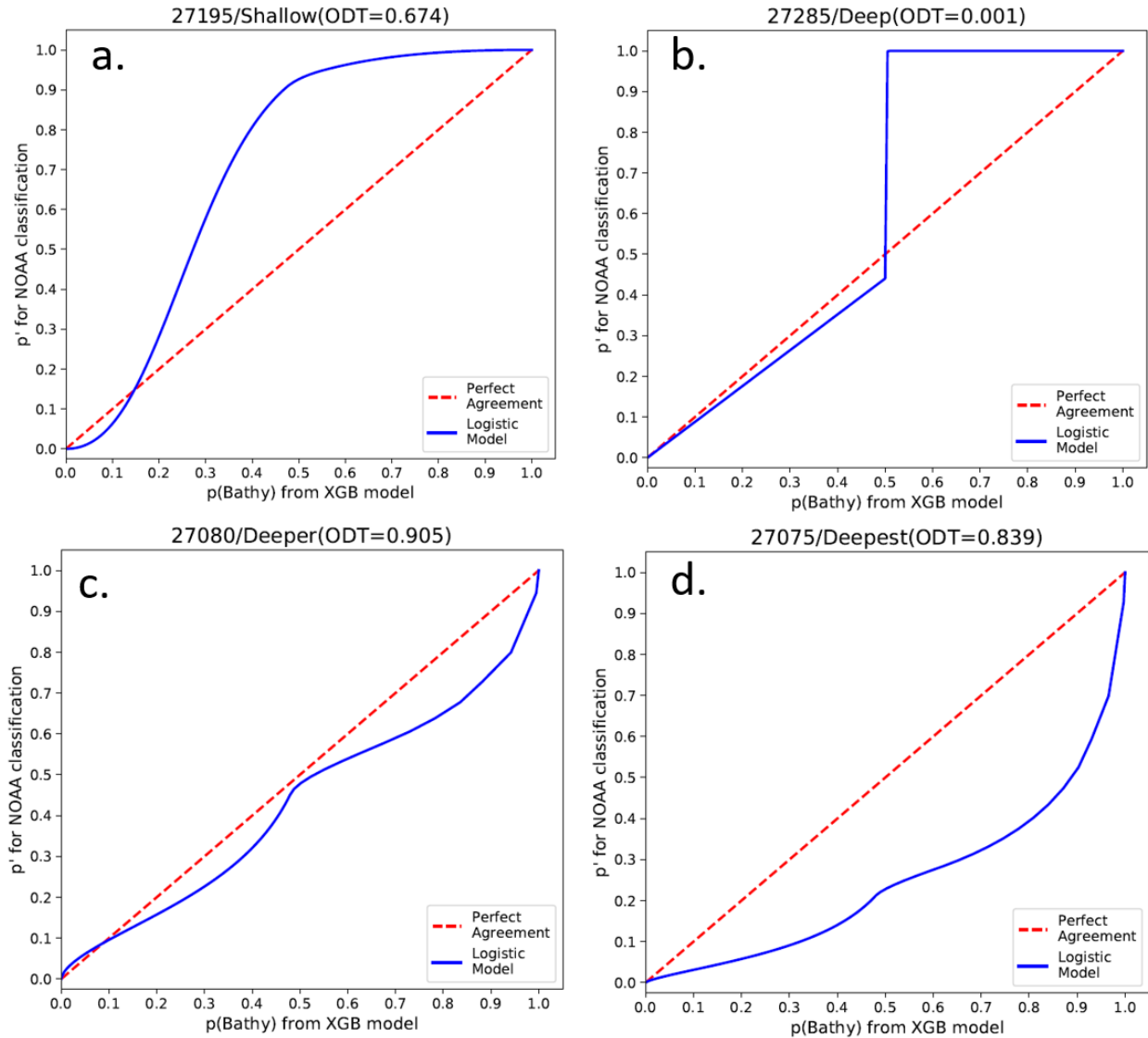
439   some tiles (Figure 6).

440
441 Figure 6. Agreement between p(Bathy) from NOAA classification and CHRT-based XGB model using logistic
442 agreement models. ("ODT" is the tile-specific optimal decision threshold.)

443 The graphs of *p'* vs. *p(Bathy)* suggest bias along the range of *p(Bathy)* values whose magnitude varies by

444 tile.  For Tiles 27080/Deeper (Fig. 6c) and 27075/Deepest (Fig. 6d), relative to the NOAA *Bathy/NotBathy*

445 classification XGB *p(Bathy)* values are overestimated over the entire range, resulting in a relatively large

446 number of false positives (FPs).  Given that these are the two tiles having depths that exceed lidar's

447 penetration capability, we hypothesize that the FPs are spatially concentrated on the deeper edges of

448 areas that NOAA identified as bathymetry; this will be examined explicitly.  The XGB model for Tile

449 27285/Deep (Fig. 6b) performs reasonably well below the ODT but severely underestimates *p(Bathy)*

450 above the ODT.  For practical purposes, this may not be problematic.  This indicates that, according to the

451 XGB model, any sounding whose *p(Bathy)* is above the ODT is 'definitely' *Bathy*. Accordingly, all pulse

20

452    returns having a *p(Bathy)* value greater than the ODT will be classified as *Bathy* – which is likely to be the

453    correct classification for the vast majority of such soundings.  Tile 27195/Shallow shows that below a

454    *p(Bathy)* value of about 0.15 the XGB model overestimates *p(Bathy)* thereby producing FP errors, but

455    higher *p(Bathy)* values are underestimates thereby resulting in FN errors.  That Fig. 5b does not indicate

456    a large number of FNs for Tile 27195/Shallow suggests that examination of the spatial distribution of FNs

457    (and FPs) might be particularly useful for this tile.

458    Examination of the geographic distribution of the differences between the ML-based and reference

459    classifications is the second method of characterizing misclassification errors.  To examine the spatial

460    distribution of the FNs and FPs, each tile was divided into 20 m pixels.  If there is no spatial bias, the errors

461    will be distributed across each tile as the lidar pulse returns are – i.e., areas having a high density of pulse

462    returns should have a comparably high density of FNs and FPs. To determine if the densities of pulse

463    returns and errors were similar and therefore not spatially biased, the differences between the percent

464    of total lidar pulse returns and percent of FNs and FPs in each pixel can be calculated with negative values

465    indicating an 'excess' of FNs or FPs.  These differences can then be displayed spatially (Figures 7, 8, 9, and

466    10) such that negative/brown values indicate 'too many' FNs or FPs, while positive/green values represent

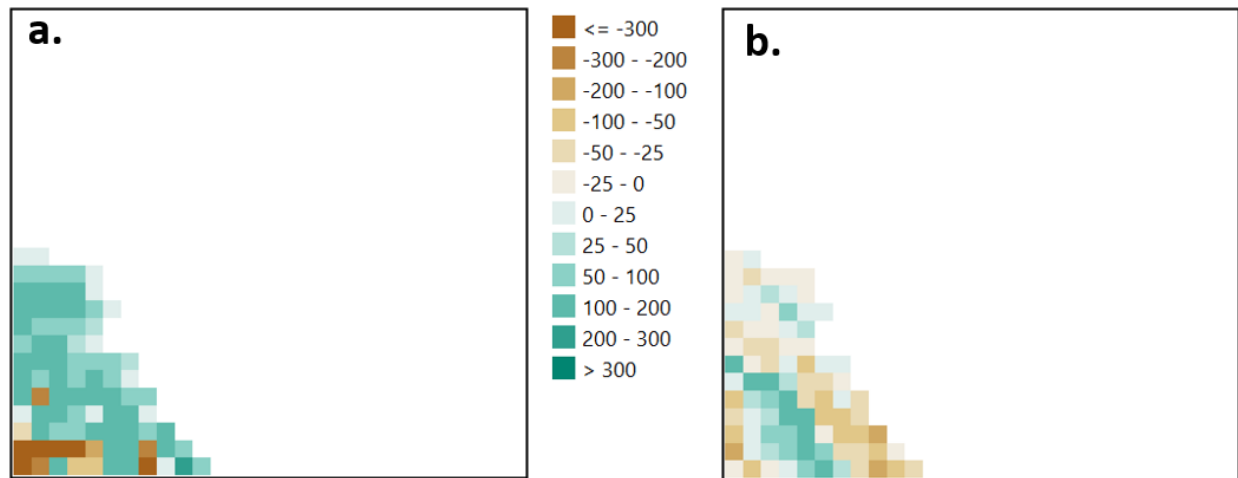467    'too few.'



468
469    Figure 7. Tile 27195/Shallow. a. Difference between percent of *Bathy* points and False Negatives (FNs) in
470    each pixel (times 100).  b. Difference between percent of *NotBathy* points and False Positives (FPs) in each
471    pixel (times 100). (Negative values indicate an "excess" of FNs or FPs.)

472    The pattern for FNs (undetected *Bathy*) for Tile 27195/Shallow (Fig. 7a) is unexpected: *Bathy* is fairly

473    accurately detected in the shallower northeast edge of data (relatively few FNs) where there are also

474    about the expected number of FPs (Fig. 7b), but there are 'too many' FNs in the deeper southwestern

475    portion.  Also of interest is that there is an area (the green northwest-to-southeast band in Fig. 7b)

476   between the shallow northeast and deeper southwest where there are 'too few' FPs.  It is also notable

477   that the magnitude of differences for FPs is less than for FNs as indicated by the more muted colors in Fig.

478   7b.

479   The information for the other tiles can be interpreted similarly:

480   ***Tile 27285/Deep (Figure 8)***: There is an 'excess' of FNs (undetected *Bathy*) in the southeastern area (Fig.

481   8a) which is the shallowest area of the tile, and an excess of FPs (erroneously detected *Bathy*) in the

482   northwest.  (The dark green north-south bands in Fig. 8a correspond to areas of flight line overlap where

483   sounding density is abnormally high.)  The magnitude of differences is greater for FNs than for FPs as

484   indicated by the more muted colors for the latter (Fig. 8b).  Noting that in practical terms FNs are

485   potentially more serious than FPs, the southeast of this tile may be an area where continuous

486   improvement efforts should be concentrated, although verifying the FPs in the northwest would also lead

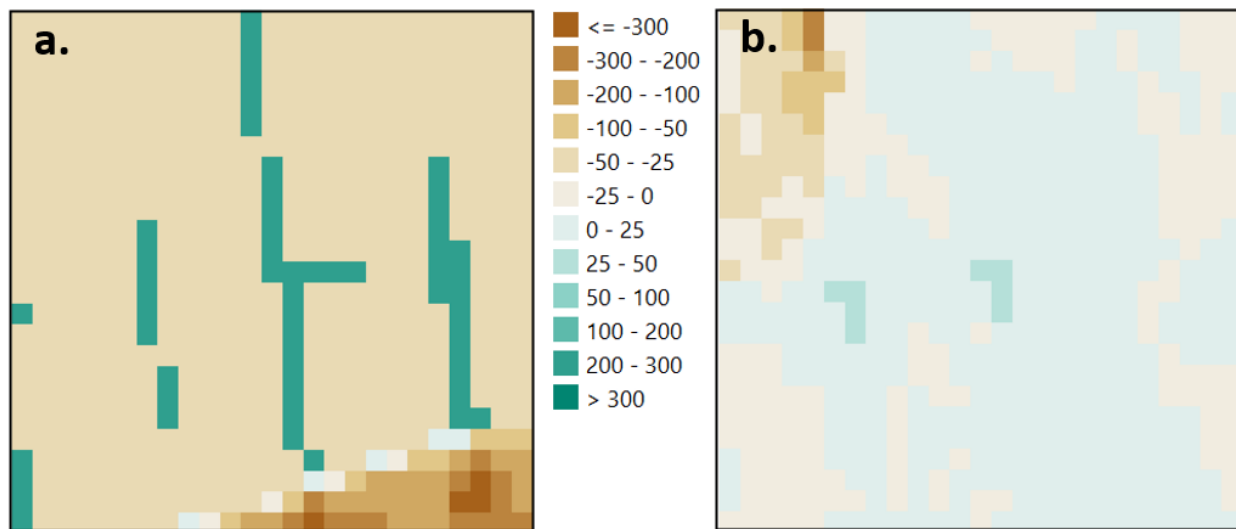487   to better accuracy for that portion of Tile 27285/Deep.



488
489   Figure 8. Tile 27285/Deep. a. Difference between percent of *Bathy* points and False Negatives (FNs) in
490   each pixel (times 100).  b. Difference between percent of *NotBathy* points and False Positives (FPs) in each
491   pixel (times 100). (Negative values indicate an "excess" of FNs or FPs.)

492   ***Tile 27080/Deeper (Figure 9)***: There is an 'excess' of FNs in the northeastern and southwestern quadrants

493   of this tile (Fig. 9a).  These are the shallowest area of this tile (see Fig. 1c) and are generally surrounded

494   by areas of relatively low differences (muted greens and browns).  Figure 9b indicates that 'too many' FPs

495   are present on the eastern edge of this tile suggesting that either the XGB classification identifies too

496   many *Bathy* soundings in this area, that the NOAA classification does not identify enough, or both.  The

497   'too few' FNs on the northwestern edge of *Bathy* pulse returns (Fig. 1c) coupled with the low level of FPs

498    in this area suggest that for this tile, the overall error rate is low at the limits of lidar light penetration.

499    The 'too high' number of FPs on the shallowest eastern edge of this tile may suggest that NOAA's

500    classification procedures, the XGB model, or both could be improved by further studying this area and the

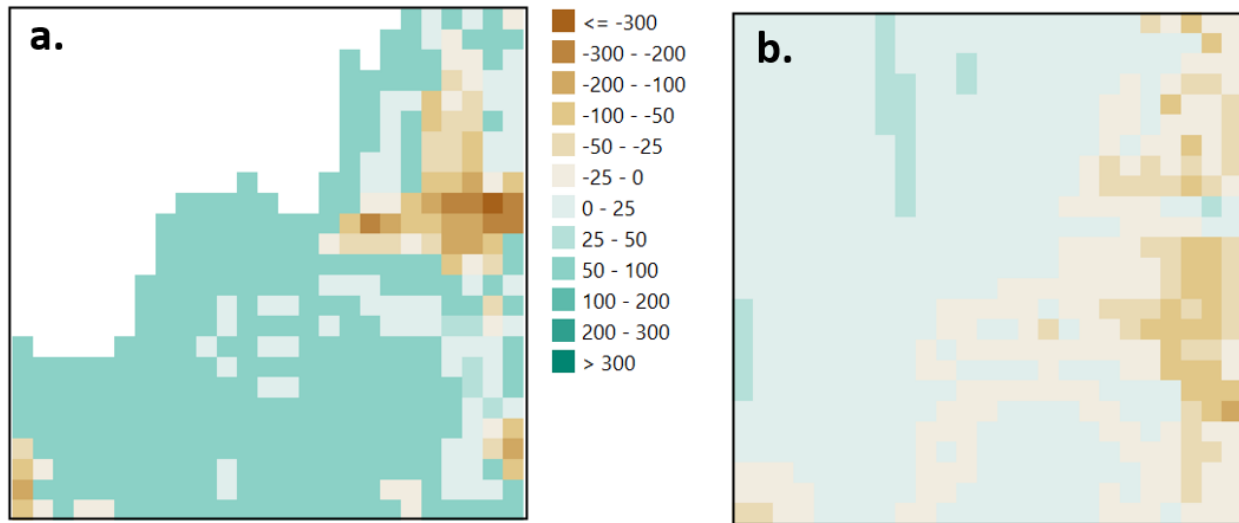501    depths it represents.



502
503    Figure 9. Tile 27280/Deeper. a. Difference between percent of *Bathy* points and False Negatives (FNs) in
504    each pixel (times 100).  b. Difference between percent of *NotBathy* points and False Positives (FPs) in each
505    pixel (times 100). (Negative values indicate an "excess" of FNs or FPs.)
506

507    **Tile 27075/Deepest (Figure 10)**: There is an 'excess' of FNs – undetected *Bathy* -- in the northeastern area

508    (Fig. 10a) which is the shallowest area and the area with the highest density of *Bathy* lidar pulse returns

509    (see Figure 1d).  There are 'too many' FPs in the same area (Fig. 10b) further suggesting that *Bathy* may

510    be under-detected in this area (and/or that the XGB model performs poorly in this area).  Interestingly,

511    however, FPs are fairly widely distributed spatially.  This distribution of FPs strongly suggests that

512    undetected *Bathy* soundings might be present throughout the tile.  It is possible that the widely

513    distributed FPs result from NOAA adopting a conservative approach to extracting bathymetry from areas

514    near the limit of lidar ocean penetration.  Regardless of the reason, these results demonstrate how this

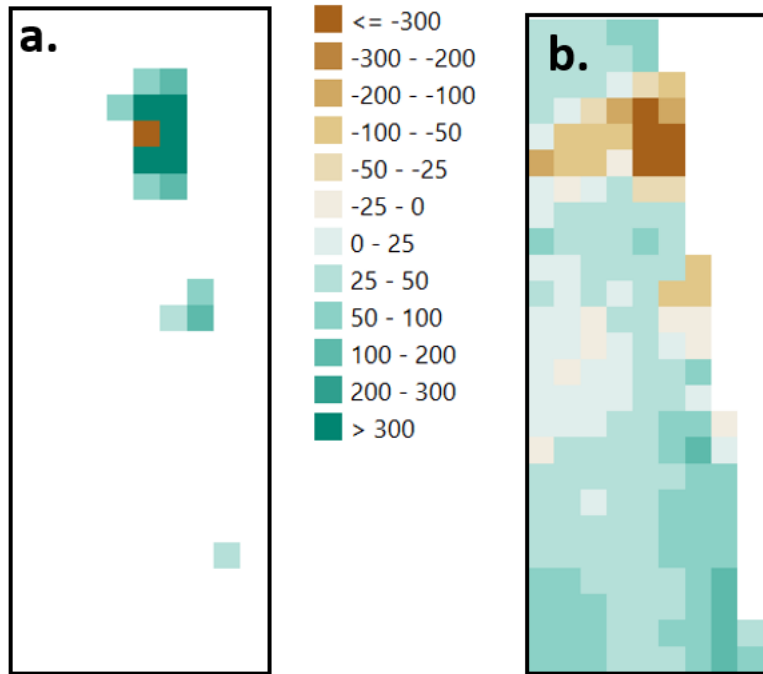515    type of analysis could help to guide continuous improvement.

Figure 10. Tile 27075/Deepest. a. Difference between percent of *Bathy* points and False Negatives (FNs) in each pixel (times 100).  b. Difference between percent of *NotBathy* points and False Positives (FPs) in each pixel (times 100). (Negative values indicate an "excess" of FNs or FPs.)

## 5.  Summary and Conclusions

Lidar soundings that identify bathymetry could be extracted from lidar point clouds without the need for an *a priori* estimate of depth with average global accuracies, true positive rates, and true negative rates of 93% compared to a reference classification for four 500 m-by-500m lidar data tiles located near Key West, Florida.   These 'accuracies' are achieved relative to a reference classification that is used operationally but that also has an unknown level of uncertainty.   The accuracy of a preliminary *Bathy/NotBathy* classification derived solely from a density-based algorithm coupled with unsupervised clustering was improved for three of the tiles by fitting and applying a machine learning extreme gradient boosting model to produce a final *Bathy/NotBathy* classification.  Models for each tile were fit using the preliminary classification as the dependent variable and 14 SAD variables such as the intensity and incidence angle of each pulse return as independent variables.  Though depth was consistently the most important SAD variable, the models for all four tiles contained at least 11 SAD variables indicating that the information that distinguishes between *Bathy* and *NotBathy* soundings is dispersed among numerous SAD variables with no consistency across all tiles.  Moreover, the information that distinguishes between *Bathy* and *NotBathy* soundings is spread among SAD that quantify pulse reflectance characteristics and airplane stability with the importance of individual SAD.

24

537    Two methods were employed to characterize differences between the reference and the ML-based

538    classification in feature/statistical and geographic space.   These are exemplified by feature-space

539    information in Table 7 and Figure 6, and spatial information presented in Figures 7 through 10.  One use

540    of this information would be refining the XGB final classification.  However, a potentially more valuable

541    application would be continuous improvement of NOAA processing methodology.  Because the true level

542    of accuracy in the XGB and NOAA classifications is unknown, the results in Table 7, Figure 6, and Figures 7

543    through 10 may be viewed as identifying differences between two independent classifications rather than

544    differences against 'truth.'  The differences might be due to weaknesses in either or both classification(s),

545    or the large-difference areas may be where bathymetry is simply difficult to extract.  Knowing this could

546    lead to revisions in XGB and/or NOAA classification procedures depending on confidence in a

547    classification, the severity or potential practical consequences of the observed differences, and a variety

548    of other factors.

549    In closing, the dual objectives of this work are recalled: to diminish the impacts of extracting bathymetry

550    from lidar sounding clouds for shallow water using machine learning, and to achieve this without the

551    circularity of needing a pre-existing classification or even depth estimate.  If one accepts that the NOAA

552    classification used for evaluation is an authoritative – but not error-free – reference, we argue that these

553    objectives have been achieved across a range of depth and data conditions.  While refinement of methods

554    and a better understanding of the nature of errors can certainly provide improvements, this work at least

555    provides a workflow to decrease time-consuming manual effort in extracting bathymetry from lidar

556    sounding clouds.

557

558    **List of Abbreviations**

559    CHRT – CUBE (Calder and Mayer 2003) with Hierarchical Resolution Technique (Calder and Rice 2017)

560    EN – Estimation Node for the density-based algorithm

561    FN, FNR – False Negative Rate (*Bathy* soundings erroneously identified as *NotBathy*)

562    FP, FPR – False Positive Rate (*NotBathy* soundings erroneously identified as *Bathy)*

563    MD – Mahalanobis Distance (used for outlier analysis)

564    ML – Machine Learning

565    MLD – Most Likely Depth

566    NOAA – National Oceanic and Atmospheric Administration

567    PDT – Probability Decision Threshold

568    SAD – Sounding Attribute Data

569     TN – True Negative Rate (*NotBathy* soundings correctly identified as *NotBathy*)

570     TP – True Positive Rate (*Bathy* soundings correctly identified as *Bathy*)

571     UTM – Universal Transverse Mercator

572

576     **Bibliography**

577     Agrifiotis, P., D., Skarlatos, A., Georgopoulos, and K., Karantzalos. 2019a. Depthlearn: learning to correct
578         the refraction on point clouds derived from aerial imagery for accurate dense shallow water
579         bathymetry based on SVMs-fusion with LiDAR point clouds. *Remote Sensing* 11: 1225, 31 pp. DOI:
580         10.3390/rs11192225.

581     Agrifiotis, P., D., Skarlatos, A., Georgopoulos, and K. Karantzalos. 2019b. Shallow water bathymetry
582         mapping from UAV imagery based on machine learning. *The International Archives of the*
583         *Photogrammetry, Remote Sensing and Spatial Information Sciences*, XLII-2/W10, pp. 9–16. DOI:
584         https://doi.org/10.5194/isprs-archives-XLII-2-W10-9-2019.

585     American Society for Photogrammetry and Remote Sensing. 2013. LAS Specification Version 1.3-R13 (15
586         July 2013).

587     Andersen, M., A. Gergely, Z. Al-Hamdani, F. Steinbacher, L. Larsen, and V., Ernstsen. 2017. Processing and
588         performance of topobathymetric lidar data for geomorphometric and morphological classification in
589         a high-energy tidal environment. *Hydrology and Earth System Sciences* 21: 43-63. DOI: 10.5194/hess-
590         21-43-2017.

591     Birkeback, M., F. Eren, S. Pe'eri, N. Weston. 2018. The effect of surface waves on airborne lidar bathymetry
592         (ALB) measurement uncertainties. *Remote Sensing* 10: 453, 19 pp. DOI: 10.3390/rs10030453.

593     Brzank, A., C. Heipke, J., Goepfert, and U. Soergel. 2008. Aspects of generating precise digital terrain
594         models in the Wadden Sea from lidar-water classification and structure line extraction. *ISPRS Journal*
595         *of Photogrammetry and Remote Sensing* 63: 510-528.

596     Calder, B., and L. Mayer. 2003. Automatic processing of high-rate, high-density multibeam echosounder
597         data. *Geochemistry, Geophysics, Geosystems*, 4(6), 1048 (22 pp.). DOI: 10.1029/2002GC000486.

598     Calder, B., and G. Rice. 2017. Computationally efficient variable resolution depth estimation. *Computers*
599         *& Geosciences* 106: 49-59. DOI: dx.doi.org/10.1016/j.cageeo.2017.05.013.

600     Calvert, J., J. String, M., Service, C., McGonigle, and R. Quinn. 2015. An evaluation of supervised and
601         unsupervised classification techniques for marine benthic habitat mapping using multibeam
602         echosounder data, *ICES Journal of Marine Science* 72(5): 1498-1513. DOI:
603         doi.org/10.1093/icesjms/fsu223.

604     Collin, A., P. Archambault, and B., Long. 2008. Mapping shallow water seabed habitat with the SHOALS.
605         *IEEE Transactions on Geoscience and Remote Sensing* 46(10): 2947-2955. DOI:
606         10.1109/TGRS.2008.920020.

607     Dietrich, J., 2017. Bathymetric structure-from-motion: extracting shallow stream bathymetry from multi-
608         view stereo photogrammetry. *Earth Surface Processes and Landforms* 42(2): 355-364.

609     Eren, F., S. Pe'eri, Y. Rzhanov, and L. Ward. 2018. Bottom characterization by using airborne lidar
610        bathymetry (ALB) waveform feature obtained from bottom return residual analysis. *Remote Sensing*
611        *of Environment* 206: 260-274. DOI: doi.org/10.1016/j.rse.2017.12.035.

612     Fernandez-Diaz, F., C. Glennie, W. Carter, R. Shrestha, M. Sartori, A. Singhania, C. Legleiter, and B.
613        Overstreet. 2014. Early results of simultaneous terrain and shallow water bathymetry mapping using
614        a single-wavelength airborne LiDAR sensor. *IEEE Journal of Selected Topics in Applied Earth*
615        *Observations and Remote Sensing* 7(2): 623-635. DOI: 10.1109/JSTARS.2013.2265255.

616     Friedman, J. 2001. Greedy function approximation: a gradient boosting machine. *Annals of Statistics* 29(5):
617        1189–1232.

618     Gholamalifard, M., T. Kutser, A. Esmaili-Sari, A., Abkar, and B. Naimi. 2013. Remotely sensed empirical
619        modelling of bathymetry in the Southeastern Caspian Sea. *Remote Sensing* 5(6): 2746-2762. DOI:
620        10.3390/rs5062746.

621     Heritage, G., and D. HetheringtonD. 2007. Towards a protocol for laser scanning in fluvial geomorphology.
622        *Earth Surface Processes and* Landforms 32(1): 66-74.

623     Hickman, G., and J. Hogg. 1969. Application of an airborne pulsed laser for near shore bathymetric
624        measurements. *Remote Sensing of Environment* 1: 47-58.

625     Höfle B., and M. Rutzinger. 2011. Topographic airborne LiDAR in geomorphology: a technological
626        perspective. *Zeitschrift für Geomorphologie* 55(2): 1-29 (in English).

627     Jawak, S., S. Vadlamani, and A. Luis. 2015. A synoptic review on deriving bathymetry information using
628        remote sensing technologies: models, methods, comparisons. *Advances in Remote Sensing* 4(2):147-
629        162. DOI: 10.4236/ars.2015.42013.

630     Kashani, A., M. Olsen, C. Parrish, and N. Wilson. 2015. A review of LIDAR radiometric processing: from *Ad*
631        *Hoc* intensity correction to rigorous radiometric calibration *Sensors* 15(11): 28099-28128. DOI:
632        doi.org/10.3390/s151128099.

633     Kerr, J., and S. Purkis. 2018. An algorithm for optically deriving water depth from multispectral imagery
634        in coral reef landscapes in the absence of ground-truth data. *Remote Sensing of Environment* 210:
635        307-324. DOI: doi.org/10.1016/j.rse.2018.03.024.

636     Kinzel, P., C. Legleiter, and P. Grams. 2021. Field evaluation of a compact, polarizing topo-bathymetric
637        lidar across a range of river conditions. *River Research and Applications* 13 pp. DOI:
638        doi.org/10.1002/rra.3771.

639     Kinzel, P., C. Legleiter, and J. Nelson. 2013. Mapping river bathymetry with a small footprint green LiDAR:
640        applications and challenges. *Journal of the American Water Resources Association* (JAWRA) 49(1):
641        183-204. DOI: 10.1111/jawr.12008.

642     Kogut, T., and M. Weistock. 2019. Classifying airborne bathymetry data using the Random Forest
643        algorithm. *Remote Sensing Letters* 10(9): 874-882. DOI: doi.org/10.1080/2150704x.2019.1629710.

644     Kutser, T., J. Hedley, C. Giardino, C., Roelfsema, and V. Brando. 2020. Remote sensing of shallow waters –
645        a 50-year retrospective and future directions. *Remote Sensing of Environment* 240: 111619, 18 pp.
646        DOI: doi.org/10.106/j.rse.2019.111619.

647     Lecours, V., M. Dolan, A., Micallef, and V. Lucieer. 2016. A review of marine geomorphometry, the
648        quantitative study of the sea floor. *Hydrology and Earth System Sciences* 20: 3207-3244, DOI:
649        https://doi.org/10.5194/hess-20-3207-2016.

650    Li, J., D. Knapp, S., Schill, C., Roelfsema, S., Phinn, M., Silman, J., Mascaro, and G. Asner. 2019. Adaptive
651        bathymetry estimation for shallow coastal waters using Planet Dove satellite. *Remote Sensing of*
652        *Environment* 232: 111302 (14 pp.). DOI: doi.org/10.1016/j.rse.2019.111302

653    Liu, S., Y. Gao, W. Zheng, and X. Li. 2015. Performance of two neural network models in bathymetry.
654        *Remote Sensing Letters* 6(4): 321–330. DOI:10.1080/2150704X.2015.1034885.

655    Lowell, K., B. Calder, and A. Lyons. 2021. Measuring shallow-water bathymetric signal strength in lidar
656        point attribute data using machine learning. *International Journal of Geographical Information*
657        *Science* (in press). Published on-line. DOI:10.1080/13658816.2020.1867147.
658

659    Lyzenga, D., N. Malinas, and F. Tanis. 2006. Multispectral bathymetry using a simple physically based
660        algorithm. *IEEE Transactions on Geoscience and Remote Sensing* 44(8): 2251-2259.

661    Maas, H.-G., D. Mader, K., Richter, and P. Westfeld. 2019. Improvements in lidar bathymetry data analysis.
662        *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*,
663        XLII-2/W10: 113-117–16. DOI: doi.org/10.5194/isprs-archives-XLII-2-W10-113-2019.

664    McFadden, D. 1974. Conditional logit analysis of qualitative choice behavior. In *Frontiers in Econometrics*,
665        Academic Press, (P. Zarembka, ed.), pp. 105-142.

666    McQueen, J. 1967. Some Methods for classification and Analysis of Multivariate Observations.
667        *Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability.* University of
668        California Press. pp. 281–297.

669    Mahalanobis, P. 1936. On the generalized distance in statistics.  *Proceedings of the National Institute of*
670        *Sciences of India* 2(1): 49–55.

671    Mandlburger, G., Pfennigbauer, M., Schwarz, R., Flory, S., and L. Nussbaumer. 2020. Concept and
672        performance evaluation of a novel UAV-borne topo-bathymetric LiDAR sensor.  *Remote Sensing*
673        12(6):986. DOI: doi.org/10.3390/rs12060986.

674    Misra, A., Z. Vojinovic, B., Ramakrishnan, A., Luijendijk, and R. Ranasinghe. 2018. Shallow water
675        bathymetry mapping using support vector machine (SVM) technique and multispectral imagery.
676        *International Journal of Remote Sensing* 39(13): 4431-4450. DOI:
677        doi.org/10.1080/01431161.2017.1421796.

678    Mitchell, S., and J. Thayer. 2014. Ranging through shallow semitransparent media with polarization lidar.
679        *Journal of Atmospheric Ocean Technology* 31: 681-697. DOI: doi.org/10.1175/jtech-d-13-00014.1.

680    Nagle, D., and C. Wright. 2016. *Algorithms used in the Airborne Lidar Processing System (ALPS)*. United
681        States Dept. of the Interior/ United States Geological Survey, Open File Report 2016-1046, 45 pp.

682    Niroumand-Jadidi, M., A. Vitti, and D. Lyzenga. 2018. Multiple optimal depth predictors analysis (MODPA)
683        for river bathymetry: findings from spectroradiometry, simulations, and satellite imagery. *Remote*
684        *Sensing of Environment* 218: 132-147. DOI: //doi.org/10.1016/j.rse.2018.09.022.

685    Okhrimenko, M., and C. Hopkinson. 2020. A simplified end-user approach to lidar very shallow water
686        bathymetric correction. *IEEE Geoscience and Remote Sensing Letters* 17(1): 3-7.

687    Pacheco, A., J. Horta, C., Loureiro, and O. Ferreira. 2015. Retrieval of nearshore bathymetry from Landsat
688        8 images: a tool for coastal monitoring in shallow waters. *Remote Sensing of Environment* 159:102–
689        116. DOI: 10.1016/j.rse.2014.12.004.

690    Pe'eri, S., and W. Philpot. 2007. Increasing the existence of very shallow LIDAR measurements using the
691        red-channel waveforms. *IEEE Transactions on Geoscience and Remote Sensing* 45(5): 1217-1223.

692 Pittman, S., B. Costa, and T. Battista. 2009. Using lidar bathymetry and boosted regression trees to predict
693    the diversity and abundance of fish and corals. *Journal of Coastal Research* 53: 27-38.

694 Schmidt, A., F. Rottensteiner, and U. Soergel. 2012. Classification of airborne laser scanning data in
695    Wadden sea areas using conditional random fields. *The International Archives of the
696    Photogrammetry, Remote Sensing and Spatial Information Sciences* XXIX-B3: 161-166. DOI:
697    https://doi.org/10.5194/isprs-archives-XLII-2-W10-9-2019.

698 Schwarz, R., G. Mandlburger, M. Pfennigbauer, and N. Pfeifer. 2019. Design and evaluation of a full-wave
699    surface and bottom-detection algorithm for LiDAR bathymetry of very shallow waters. *ISPRS Journal
700    of Photogrammetry and Remote Sensing* 150: 1-10. DOI: dx.doi.org/10.1016/j.isprsjprs.2019.02.002.

701 Steinhaus, H. 1957. Sur la division des corps matériels en parties. (English: "On the division of body
702    materials in parts.") Bulletin de l'Académie Polonaise des Sciences. (English: "Bulletin of the Polish
703    Academy of Sciences.") 4(12): 801–804. (In French.)

704 Su, D., F. Yang, Y. Ma, K. Zhang, J. Huang, and M. Wang. 2019. Classification of coral reefs in the South
705    China Sea by combining airborne LiDAR bathymetry bottom waveforms and bathymetric features.
706    *IEEE Transactions on Geoscience and Remote Sensing* 57(2): 815-828. DOI:
707    10.1109/TGRS.2018.2860931.

708 Tulldahl, H., and S. Wikström. 2012. Classification of aquatic macrovegetation and substrates with
709    airborne lidar. *Remote Sensing of Environment* 121: 347-357. DOI:
710    //doi.org/10.1016/j.rse.2012.02.004.

711 Wang, C., Q. Li, Y. Liu, G. Wu, P. Liu, and X. Ding. 2015. A comparison of waveform processing algorithms
712    for single-wavelength LiDAR bathymetry. *ISPRS Journal of Photogrammetry and Remote Sensing*.
713    101**:**22-35. DOI: dx.doi.org/10.1016/j.isprsjprs.2014.11.005.

714 Wang, L., H. Liu, H. Su, and J. Wang. 2018. Bathymetry retrieval from optical images with spatially
715    distributed support vector machines. *GIScience & Remote Sensing*. 56(3); 323-337. DOI:
716    10.1080/15481603.2018.1538620.

717 Westfeld, P., H. Maas, K. Richter, K., and R. Weiss. 2017. Analysis and correction of ocean wave pattern
718    induced systematic coordinate errors in airborne LiDAR bathymetry. *ISPRS Journal of
719    Photogrammetry and Remote Sensing* 128:314-325.

720 Xing, S., D. Wang, Q. Xu, Y. Lin, P. Li, L. Jiao, Z. Zhang, and C. Liu. 2019. A depth-adaptive waveform
721    decomposition method for airborne LiDAR bathymetry. *Sensors* 19: 5065, 28 pp. DOI:
722    10.3390/s19235065.

723 Yang, A., Z. Wu, F. Yang, D. Su, Y. Ma, D. Zhao, and C. Qi. 2020. Filtering of airborne LiDAR bathymetry
724    based on bidirectional cloth simulation. *ISPRS Journal of Photogrammetry and Remote Sensing*
725    163:49-61. DOI: doi.org/10.1016/j.isprsjprs.2020.03.004.

726 Zhao, J., X. Zhao, H. Zhang, and F. Zhou. 2017. Shallow water measurements using a single green laser
727    corrected by building a near water surface penetration model. *Remote Sensing* 9**:** 426; 18 pp.,
728    DOI:10.3390/rs9050426.