1 **A datamining approach to identifying spatial patterns of phosphorus forms**

2 **in the Stormwater Treatment Areas in the Everglades, US.**

3

4 Corstanje, R.[a]*, Grafius, D.R.[a], Zawadzka, J.[a], Moreira Barradas, J.[a], Vince, G.[b],

5 Ivanoff, D.[c], Pietro, K. [c]

6

7 *Corresponding author: roncorstanje@cranfield.ac.uk

8

9 [a]Cranfield Soil and Agrifood Institute, Cranfield University, College Road,

10 Cranfield, Bedfordshire, MK43 0AL, UK

11

12 [b]Tetra Tech Inc., 759 S. Federal Hwy., Suite 314, Stuart, FL 34994-2936, USA

13

14 [c]South Florida Water Management District, 3301 Gun Club Road, West Palm

15 Beach, FL 33406, USA

16

17 **Abstract**

18

19 The Everglades ecosystem in Florida, USA, is naturally phosphorus (P) limited,

20 and faces threats of ecosystem change and associated losses to habitat,

21 biodiversity, and ecosystem function if subjected to high inflows of P and other

22 nutrients. In addition to changes in historic hydropattern, upstream agriculture

23 (sugar cane, vegetable, citrus) and urbanization has placed the Everglades at risk

24 due to nutrient-rich runoff. In response to this threat, the Stormwater Treatment

25 Areas (STAs) were constructed along the northern boundary of the Everglades as

engineered ecological systems designed to retain P from water flowing into the Everglades. This research investigated data collected over a period from 2002 to 2014 from the interior of the STAs using data mining and analysis techniques including a) exploratory methods such as Principal Component Analysis to test for patterns and groupings in the data, and b) modelling approaches to test for predictive relationships between environmental variables. The purpose of this research was to reveal and compare spatial trends and relationships between environmental variables across the various treatment cells, flow-ways, and STAs. Common spatial patterns and their drivers indicated that the flow-ways do not function along simple linear gradients; instead forming zonal patterns of P distribution that may increasingly align with the predominant flow path over time. Findings also indicate that the primary drivers of the spatial distribution of P in many of these systems relate to soil characteristics. The results suggest that coupled cycles may be a key component of these systems; i.e. the movement and transformation of P is coupled to that of nitrogen (N).

**Keywords**: phosphorus, data mining, stormwater treatment areas, constructed wetland, Everglades, water quality

**1. Introduction**

The Stormwater Treatment Areas (STAs), located around the northern boundary of the Everglades in Florida, USA, were constructed over a period from 1994 to 2013. As a set of engineered ecological systems, the general purpose and function of the STAs is to reduce phosphorus (P) in runoff water prior to

51    discharging to the Everglades Protection Area. They consist of a series of

52    shallow, freshwater marshes divided into flow-ways and treatment cells by

53    interior levees and control structures, populated with emergent or submergent

54    aquatic vegetation (EAV and SAV, respectively) (Chen et al., 2015). The

55    Everglades as a system is naturally P limited (Entry, 2014; McCormick et al.,

56    1996), and so the water it receives must meet stringent requirements for ultra-

57    low levels of water P (Pietro and Ivanoff, 2015). Since 1995, the STAs have

58    treated approximately 16.5 billion $m^3$ inflow volume, retained approximately

59    1,727 metric tons (mt) of total phosphorus (TP), lowering phosphorus surface

60    water concentrations from an overall annual TP of 140 micrograms per liter (μg

61    $L^{-1}$) to 37 μg $L^{-1}$ (flow weighted mean; South Florida Water Management District,

62    2015), and improving further in most recent years to exhibit outflow

63    concentrations averaging between 15-25 μg $L^{-1}$ (South Florida Water

64    Management District et al., 2015). STA-2 and STA-3/4 are two of the best

65    performing STAs, and have recorded reductions in surface water P from 100 and

66    87 μg $L^{-1}$ at inflow structures, respectively, to 23 and 18 μg $L^{-1}$ at outflow (Pietro

67    and Ivanoff, 2015).

68

69    The STAs are wetland systems, and the controls on the P removal process are

70    therefore set by the internal biogeochemical, ecological and physical processes

71    and conditions in each cell, in each STA (Ivanoff et al., 2013). Phosphorus

72    reduction from each STA must be maximized in order to meet stringent

73    regulatory effluent limits, which implies that these natural processes must be

74    manipulated (engineered) to maximize P retention. Phosphorus in surface water

75    can have various forms; from soluble reactive to forms of organic and particulate

76  P with varied degrees of recalcitrance (Reddy and DeLaune, 2008). The retention

77  of P in these systems needs to therefore consider these different forms.

78

79  There are abiotic processes of P retention, including P sorption to the STA soil

80  particulates (Reddy et al., 1999) and particulate (co)-precipitation with cations

81  such as calcium (Ca), magnesium (Mg), iron (Fe), and aluminium (Al) (Malecki-

82  Brown et al., 2007). Factors that influence these processes are surface flow rate

83  and path (Kadlec and Wallace, 2009) but also water and soil chemistry (e.g.

84  concentrations of Ca, Mg, Fe and Al), pH, and the oxidation reduction potential

85  (Reddy et al., 1999).  Ideally this P then gets buried, or retained by the sediment

86  within the wetland, resulting in gradually lower soil P-levels as water flows from

87  the inflow point towards the outflow points (P gradient), similar to what has

88  been observed in the nearby Water Conservation Area 2A (DeBusk et al., 1994).

89  There are circumstances under which P is transported along the hydrologic

90  gradient due to sediment re-suspension, P desorption from the sediment matrix,

91  or poor vegetation condition. In properly performing STAs, these are limited and

92  water column P could be reduced further down the flow-way, reducing the slope

93  of the gradient.  Uptake and retention of P by plants is generally (though not

94  exhaustively; dependent upon plant type) considered to be short-term and rapid;

95  while abiotic/physical retention processes tend to be longer term and are

96  considered to account for 50-70% of permanent storage (Richardson, 1999).

97

98  Biological cycling of P involves direct uptake of available P by plant and

99  microbial communities (Newman et al., 2001) to meet their physiological

100 requirements, action of extracellular enzymes on complex organic P to release P

101     uptake (Corstanje et al., 2007) and the release of P from the biological

102     decomposition of organic material. Under anaerobic environments,

103     decomposition of organic material is slow, resulting in formation and accretion

104     of peat; forming another sink for P as long as the peat remains intact. Biological P

105     cycling and the resulting spatial distribution of the different forms of P is highly

106     complex, as it is driven by coupled P, N and C cycles; determined by redox

107     conditions and characterized by the plant ecology (Chen et al., 2015; Orem et al.,

108     2014; Reddy et al., 2011).

109

110     Extensive sampling has been conducted over a period from 2002 to 2014, in

111     which soil, surface water and macrophytes have been sampled within the STA

112     cells, resulting in a large dataset of observations. Coupled with hyper-spectral

113     measurements made through various aerial surveys, the results comprise a fairly

114     comprehensive dataset on the spatial variation in key components of the STA

115     ecosystem. Here, we report on a broad scale analysis of these datasets, in order

116     to determine common trends across the various flow-ways in the STAs, and in

117     individual STAs. The expectation here is that common biogeochemical processes

118     will generate common multivariate patterns across STAs. We then considered,

119     given the extent and comprehensiveness of the datasets under consideration,

120     implications for future monitoring of these systems.

121

122     **2. Materials and Methods**

123     **2.1. Study Area**

124

125    The STAs, operated by the South Florida Water Management District, cover an

126    effective treatment area of circa 230 km². There are five STAs: STA-1E, STA-1W,

127    STA-2, STA-3/4, and STA-5/6 (Figure 1); STA-5/6 was formerly two separate

128    STAs until water year (WY) 2010. The STAs vary in size and location, and each is

129    constructed with sets of interconnected cells forming treatment 'flow-ways'.

130    Data from surface water (sampled along internal transects within the treatment

131    cells), floc (i.e. flocculant; loosely clumped particles either suspended in the

132    water column or resting atop the soil, analogous to litter in terrestrial systems),

133    and soil collected within the various cells were available for analysis, and have

134    been previously described and used to evaluate conditions within the STAs (e.g.

135    Pietro and Ivanoff, 2015; Reddy et al., 2009). Normalized Difference Vegetation

136    Index (NDVI) and vegetation class and habitat maps were derived from recent-

137    year hyper-spectral imagery at a resolution of approximately 1 square foot to

138    represent the approximate current state of vegetation within the cells. The

139    available datasets were diverse in spatial extents, subjects (e.g. soil samples,

140    surface water transects, vegetation coverage) and data types (e.g. categorical vs.

141    continuous), necessitating a data mining approach capable of addressing this

142    diversity. Below we describe the structure of each STA; specifics of data

143    availability are described in the sections that follow.

144

145    STA-1E began full operation in 2006-2007 and consists of three flow-ways;

146    Eastern, Western, and Central. Due to data availability only the Central Flow-way

147    was analyzed here. STA-1W's Eastern and Western flow-ways were in operation

148    from 1994 as the Everglades Nutrient Removal (ENR) project, with an additional

149    Northern flow-way constructed in 2000. All three flow-ways were analyzed. STA-

150  2 Cells 1-3, each single-cell flow-ways, were operational from 2000 onwards.

151  Additional cells, 4-8, involve multi-cell flow-ways and became operational

152  between 2008 and 2012 but were not studied here due to insufficient data

153  availability. STA-3/4 consists of three flow-ways (Flow-ways 1, 2 and 3) and

154  became operational in 2004; all were included in analysis. STA-5 originally

155  consisted of three flow-ways, denoted Flow-ways 1, 2 and 3; each consisting of a

156  combination of two cells. Flow-ways 1 and 2 became operational in 1999; Flow-

157  way 3 in 2008. Flow-ways 4 and 5 were later added, flow-capable in 2010, but

158  not studied here. Combination with STA-6 to form STA-5/6 added three

159  additional flow-ways; 6, 7 and 8, of which Flow-ways 7 and 8 are single cell flow-

160  ways (operational in 1998), and Flow-way 6 (not analyzed) couples two cells (6-

161  4, flow-capable in 2010 and 6-2, constructed in 2006).

162

163  **2.2. Data quality control**

164

165  Quality control checks were performed on all datasets at various stages of the

166  data compilation. Blank or null records were treated as no data and not zero.  For

167  soil and floc data, parameter values were reported within specific ranges of the

168  profile, typically ranging from 0 to 10 cm. Some records included data on the

169  upper profile (0-10 cm), lower profile (10-30 cm), and full profile combined (0-

170  30 cm). In some cases soil nutrients within selected STA cells were measured at

171  variable depth increments (e.g. 0-2, 2-4, 4-6 cm, etc.). In such cases, all

172  parameters for relevant increments were averaged into a single 0-10 cm field for

173  analysis to ensure consistency across the dataset (including bulk density). In

174  some other cases, the sampling depth of the upper profile did not reach 10 cm,

175 but these were still marked as the upper profile. The full profile value was very

176 rarely given, and was calculated only for the datasets that were subsequently

177 used in the data mining analysis. In these instances, the average of the upper and

178 lower profile was used.

179

180 **2.3. Data Analysis**

181 **2.3.1. Preparation of datasets for data mining**

182

183 The following rules were applied for inclusion of the data measured within the

184 STAs: (1) There must be at least 10 observations for a given STA cell and year

185 (an arbitrary cutoff point but sufficient to allow the calculation of meaningful

186 statistics) and (2) There must be at least one instance of at least 10 observations

187 per year within all STA cells in a flow-way. Seasonality at temporal scales finer

188 than full years was not considered here. Additionally, any GIS data with full

189 coverage of STA cells were considered. These included vector maps of vegetation

190 class and habitat, NDVI rasters, and topography rasters representing the

191 elevation differences of the STA floor at various year intervals. The resulting

192 flow-ways included in data mining and their available data are listed in Table 1.

193 **Table 1.** List of flow-ways included in interpolation and their available data
194 including years and number of observations (n). Surface water quality data are
195 from transects internal to each treatment cell.

| STA | Flow-way | Cells | STA Data Availability |
|---|---|---|---|
| STA-1E | Central | 3 to 4N to 4S | Soil/floc (2004, 07, 09, 10; n=97) Surface water (2013; n=16) Macrophyte nutrients (2009; n=46) Hyper-spectral imagery (2011-12) |
| STA-1W | Eastern | 1A and 1B to 3 | Soil/floc (Eastern/Western FW only: 1995-97, 99; all FW: 2003-08, 10; n=1006) |
|  | Western | 2A and 2B to 4 | Surface water (2003, 04, 09-13; n=2689) |

| | | | |
|---|---|---|---|
| | Northern | 5A to 5B | Macrophyte nutrients (Eastern/Western FW only: 1996, 97; all FW: 2003, 04, 08-10; n=262) |
| | | | Hyper-spectral imagery (2011-12) |
| STA-2 | Flow-way 1 | 1 | Soil/floc (2003, 07, 09-11; n=830) |
| | Flow-way 2 | 2 | Surface water (2003-10, 13, 14; n=1126) |
| | Flow-way 3 | 3 | Macrophyte nutrients (2003, 09, 10; n=91) |
| | | | Hyper-spectral imagery (2011-12) |
| STA-3/4 | Flow-way 1 | 1A to 1B | Soil/floc (2004, 07, 10; n=1272) |
| | Flow-way 2 | 2A to 2B | Surface water (2003-10, 13, 14; n=1134) |
| | Flow-way 3 | 3A to 3B | Macropyte nutrients (2010-12; n=58) |
| | | | Hyper-spectral imagery (2011-2012) |
| STA-5/6 | Flow-way 1 | 1A to 1B | Soil/floc (FW 1/2: 2002, 03, 07-11; n=617. FW 7/8: 2003, 07-11; n=138) |
| | Flow-way 2 | 2A to 2B | Surface water (FW 1/2: 2013; n=74) |
| | Flow-way 7 | 5 | Macrophyte nutrients (FW 1/2: 2002, 03; n=147. FW 7/8: 2003; n=31) |
| | Flow-way 8 | 3 | Hyper-spectral imagery (all FW: 2011-12) |

196

197

### 2.3.2. Interpolation of flow-way data within STA cells

199

Interpolation was done using an Empirical Bayesian Kriging (EBK) algorithm. For Bayesian geostatistical analysis, we used the Gaussian Spatial Linear Mixed Model as formulated by Diggle et al. (1998) without fixed effects:

203

$$Y(s_i) = W(s_i) + \varepsilon$$

204

where the random variable $Y(s_i)$ is an $n \times 1$ vector of observed values at locations $s_1, s_1, \ldots, s_i$; $W$ represents the spatial random effect which is a Gaussian process with mean of 0, variance of $\sigma^2$ (partial sill) and correlation function $R(h; \varphi)$, for which we selected an exponential correlation function: $R(h; \varphi) = \exp(-\frac{h}{\varphi})$; and $\boldsymbol{\varepsilon}$ is an $n \times 1$ vector of errors with mean of 0 and variance of

210     $\tau^2$(nugget variance). These semivariogram parameters were estimated using

211     restricted maximum likelihood (REML). The EBK tool produced 1137 pairs of

212     interpolated and standard error maps which, together with other spatial

213     datasets available (described above in 2.3.1), were sampled with 100 randomly

214     distributed points (separated by at least fifty feet) within each STA cell.

215

216     **2.3.3. Multivariate Analysis**

217

218     Multivariate analysis used a combination of exploratory and modeling tools to

219     identify underlying patterns in the data. Within each treatment flow-way, data

220     from all available years were pooled to facilitate a single, data-rich analysis. For

221     initial calculation of summary statistics, the record set within each cell

222     containing the greatest number of observations for each year of coverage was

223     selected, and the mean and standard deviation of TP measurements were

224     calculated across all recorded years in Microsoft Excel (Microsoft, 2003). The

225     mean and standard deviation of key soil nutrients (i.e. total phosphorus, nitrogen

226     and carbon) were calculated for entire STAs. Principal components analysis

227     (PCA) and clustering analysis (CA) were used in an exploratory mode using JMP

228     (SAS, 2013); PCA to determine the main axis of variation the datasets, and CA to

229     determine if there were any meaningful groups in the observations. The primary

230     goals were: (a) to determine if there are any consistent main drivers of variation

231     across the flow-ways (i.e. do the flow-ways and STAs behave consistently across

232     the board, or is each a unique system responding to unique operational

233     circumstances); and (b) within each flow-way, to determine if there are natural

234     groupings of multivariate data (e.g. are observations from areas around the

235    inflow sufficiently similar in floc, soil and vegetation characteristics to cluster,

236    and sufficiently distinct from other areas). We used a combination of Ward's and

237    *k*-means clustering methods (Corstanje et al., 2009). Ward's is a minimum

238    variance, hierarchical clustering method which produces a scree plot, that in turn

239    allows us to both identify the optimal number of clusters and establish the seeds

240    which are then used to run the *k*-means clustering process. This was then

241    followed by Stepwise Canonical Discriminant (SCD) analysis in JMP (SAS, 2013)

242    to help identify the primary drivers of the clusters.

243

244    Subsequently, we applied a set of non-linear, hierarchical structured models

245    using Statistica (StatSoft, 2014) to predict surface water TP concentrations

246    (Classification and Regression Trees; CART). Where no surface water TP data

247    were available (as was the case in 10 out of 24 cells: STAs 1E, 2 Cell 2 only, 3/4,

248    and 6), floc TP was substituted as the best available indicator of TP and its

249    drivers in the flowing system. The CART approach has a number of advantages;

250    the method is not sensitive to non-normal data, it accepts categorical as well as

251    continuous data (needed as soil series and soil parent material are categorical,

252    whereas soil organic matter is continuous) and it is not confounded by the

253    presence of non-linear relationships (Breiman et al., 1984; McCune and Grace,

254    2002). Bayesian Belief Networks (BBNs), having similar advantages in their

255    ability to handle non-normal and categorical data, were also created using Netica

256    (Norsys, 2014) to predict the most recently available NDVI and TP

257    (preferentially in surface water if available, otherwise in floc or soil as described

258    above) in each cell.  BBNs are graphical probabilistic models; graphical in that

259    they represent the variables that affect the response of interest (e.g. floc or

260 surface water P) in the form of a network, and probabilistic in that the

261 relationships between the drivers and response are conditioned by a probability

262 (Taalab et al., 2015). Bayesian inference is thus based on a set of prior

263 probabilities that can be updated as new information becomes available. In this

264 case, some knowledge of potential drivers of P dynamics was available from the

265 CART analysis and a review of the existing P process literature; the network thus

266 consisted of those variables that the previous CART models identified as drivers.

267 For both CART and BBN approaches, model fitness and the strongest predictor

268 variables were of primary interest.

269

270 **3. Results**

271 **3.1. Summary Statistics**

272

273 Data on TP from internal surface water transects and TP, total carbon (TC) and

274 nitrogen (TN) from soil samples in all STAs and across all available years were

275 pooled and their summary statistics calculated (Tables 2 and 3), but

276 distributions were highly variable in terms of timing, data type, number of

277 observations, and data were not available or complete for all cells and flow-ways.

278 Cell 2A in STA-5/6 achieved the highest overall mean internal surface water TP

279 (0.216 mg L$^{-1}$) followed by STA-1W's Cell 5A (0.129 mg L$^{-1}$). The Cells with the

280 lowest mean internal surface water TP were STA-3/4's Cell 3B (0.012 mg L$^{-1}$)

281 and STA-1W's Cell 4 (0.024 mg L$^{-1}$). Variability was present in the data, both

282 within sets of records and between different years and cells; most standard

283 deviations tended to fall proportionally between 30% and 80% of their

284 associated means. Total soluble phosphorus (TSP) and soluble reactive

285 phosphorus (SRP) in internal surface water were variable in their proportional

286 relationship with TP (not shown); combined across all STAs, TSP averaged

287 roughly half of TP (59.2%) with a standard deviation of 15.4%, and SRP averaged

288 28.1% of TP with a standard deviation of 14.6%. As these statistics summarize

289 the data for entire treatment cells they do not address spatial patterns within

290 individual cells (this is explored below in section 3.3); however in flow-ways

291 composed of multiple cells, an apparent trend of decreasing mean TP was visible

292 along the length of the flow-ways from the summary statistics, evidencing the

293 removal of phosphorus from surface water as it flows through the STAs. The

294 greatest proportional drop was in Flow-way 2 in STA-5/6, where Cell 2A

295 exhibited a mean TP of 0.216 mg L$^{-1}$ and Cell 2B a mean of 0.062 mg L$^{-1}$.

296

297 **Table 1:** Summary statistics for all combined data on total surface water
298 phosphorus [mg L$^{-1}$] sampled within the STAs (internal surface water transect).
299 SD = Standard Deviation, N = number of observations. Values marked 'n/a'
300 represent cells where summary data were insufficient for calculation of
301 summary statistics.

| STA | Flow-way | Cells | Mean | SD | N |
|---|---|---|---|---|---|
| STA-1E | Central | 3 | n/a | n/a | 0 |
| | | 4N | 0.108 | 0.017 | 16 |
| | | 4S | n/a | n/a | 0 |
| STA-1W | Eastern | 1A | 0.106 | 0.049 | 8 |
| | | 1B | 0.065 | 0.044 | 159 |
| | | 3 | 0.030 | 0.018 | 95 |
| | Western | 2A | 0.123 | 0.069 | 77 |
| | | 2B | 0.047 | 0.022 | 89 |
| | | 4 | 0.024 | 0.012 | 70 |
| | Northern | 5A | 0.129 | 0.051 | 54 |
| | | 5B | 0.071 | 0.079 | 699 |
| STA-2 | Flow-way 1 | 1 | 0.044 | 0.036 | 197 |
| | Flow-way 2 | 2 | n/a | n/a | 0 |
| | Flow-way 3 | 3 | 0.034 | 0.024 | 606 |
| STA-3/4 | Flow-way 1 | 1A to 1B | n/a | n/a | 0 |
| | Flow-way 2 | 2A to 2B | n/a | n/a | 0 |
| | Flow-way 3 | 3A | 0.037 | 0.005 | 4 |

| | | 3B | 0.012 | 0.001 | 42 |
|---|---|---|---|---|---|
| STA-5/6 | Flow-way 1 | 1A | 0.064 | 0.048 | 12 |
| | | 1B | 0.031 | 0.023 | 16 |
| | Flow-way 2 | 2A | 0.216 | 0.074 | 12 |
| | | 2B | 0.062 | 0.045 | 16 |
| | Flow-way 7 | 5 | n/a | n/a | 0 |
| | Flow-way 8 | 3 | n/a | n/a | 0 |

302
303

**Table 3:** Summary statistics for all combined data on total soil phosphorus [TP; mg kg$^{-1}$], total carbon [TC; g kg$^{-1}$] and total nitrogen [TN; g kg$^{-1}$] sampled within the STAs. SD = Standard Deviation, N = number of observations.

| *STA* | *Soil TP (mg kg$^{-1}$)* | | | *Soil TC (g kg$^{-1}$)* | | | *Soil TN (g kg$^{-1}$)* | | |
|---|---|---|---|---|---|---|---|---|---|
| | *Mean* | *SD* | *N* | *Mean* | *SD* | *N* | *Mean* | *SD* | *N* |
| STA-1E | 241 | 207 | 294 | 85.2 | 63.3 | 294 | 5.7 | 4.3 | 294 |
| STA-1W | 550 | 237 | 1405 | 432 | 57.6 | 1322 | 26.5 | 4.2 | 1319 |
| STA-2 | 611 | 250 | 1166 | 392 | 51.2 | 1078 | 23.1 | 3.7 | 1078 |
| STA-3/4 | 718 | 243 | 1858 | 346 | 74.2 | 1857 | 22.1 | 5.0 | 1857 |
| STA-5/6 | 727 | 315 | 952 | 285 | 111 | 783 | 20.5 | 7.8 | 783 |

307

308

309 Data for TP, TC and TN in soil and floc across the STAs were analyzed at the STA

310 level. STA-5/6 exhibited the highest mean levels of soil TP (727 mg kg$^{-1}$), while

311 STA-1W achieved the highest values for both mean TC (432 g kg$^{-1}$) and mean TN

312 (26.5 g kg$^{-1}$). STA-1E had the lowest mean values for all three nutrients; 241 mg

313 kg$^{-1}$ TP, 85.2 g kg$^{-1}$ TC, and 5.7 g kg$^{-1}$ TN. Variability was highest in STA-5/6

314 across all three nutrients; exhibiting a standard deviation of 315 mg kg$^{-1}$ TP, 111

315 g kg$^{-1}$ TC, and 7.8 g kg$^{-1}$ TN. TP variability was lowest in STA-1E (standard

316 deviation of 207 mg kg$^{-1}$), while STA-2 displayed the lowest variability for both

317 TC (51.2 mg kg$^{-1}$) and TN (3.7 g kg$^{-1}$). Note that these statistics represent

318 averages across entire treatment cells or STAs; Table 3 reports the associated

319 variability (as standard deviations).

320

321 **3.2. Multivariate Analysis Results**

322

Principal Component Analysis (PCA) results are characteristically not straightforward to interpret and do not involve clear cutoffs to determine whether or not a component variable can be considered specifically important or unimportant, so focus was placed on determining and reporting those variables that were clearly the strongest drivers and/or recurred consistently across STAs. Results varied by cell, but the most commonly identified variables related to soil TC, soil TN, soil and floc bulk density (BD), soil and floc TP, and soil and floc ash-free dry weight (AFDW) as the greatest contributors to variability in the data (Table 4). Cluster analysis identified 3 or 4 clusters in most cells, with spatial structure to cluster membership apparent in some but not all cells (Table 5).

333

**Table 4:** Summary of the main outcomes from Principal Component Analysis (soil/floc/surface water parameters separated by semicolon). Abbreviations: total phosphorus (TP), total carbon (TC), total nitrogen (TN), bulk density (BD), sulfur (S), calcium (Ca), iron (Fe), macrophyte nutrient (macro), exchange capacity (exc), surface water (sw), alkalinity (Alk) ash-free dry weight (AFDW). Note that data availability was not consistent (e.g. few surface water observations in STA-1E Cells 3 and 4S) so PCA may not accurately reflect the importance of underrepresented variables in some cells.

| STA | Flow-way | Cell | PCA main variables | % var explained by PC1,..,PC3 |
|---|---|---|---|---|
| **STA-1E** | Central | 3 | Soil TC, TN, AFDW, BD, TP, Ca | 80.25 |
| | | 4N | Soil AFDW, BD, TC, Ca, Fe, TP | 79.58 |
| | | 4S | Soil AFDW, BD, TC, TN, TP, Ca, Fe | 77.77 |
| **STA-1W** | Northern | 5A | Soil AFDW, BD, TC, TN, TP; floc AFDW; sw TP | 83.23 |
| | | 5B | Floc BD, TC, AFDW; sw Ca, P | 68.01 |
| | Eastern | 1 | n/a* | 57.49 |
| | | 3 | Soil Al exc, Fe exc, TN, Alk, AFDW, BD, K; sw TP | 76.59 |
| | Western | 2 | Soil Fe, BD, TC; sw TP, Ca, AFDW | 81.34 |

| STA | | Cell | Driving variables | Value |
|---|---|---|---|---|
| | | 4 | Soil AFDW, Fe, TC, TN; sw Ca, TP | 81.66 |
| **STA-2** | 1 | 1 | Soil TC, TN, TP; floc TC; sw TP | 75.09 |
| | 2 | 2 | Floc BD, TC, TN, TP | 72.19 |
| | 3 | 3 | Soil macroDryWt, TC; sw Ca, TP | 77.98 |
| **STA-3/4** | 1 | 1A | Soil BD, TN; floc BD, TP | 72.54 |
| | | 1B | Soil TP, TN, BD; floc TN; sw TN | 71.78 |
| | 2 | 2A | Soil BD, TC, TN; floc BD, TC, TN | 63.12 |
| | | 2B | Soil TC, TP, TN, BD | 71.71 |
| | 3 | 3A | Soil TC, TN, BD, TP; sw Ca, TP | 74.21 |
| | | 3B | Floc TC, TN; sw Ca, P | 69.59 |
| **STA-5/6** | 1 | 1A | Soil macro TN, Fe, BD; floc DryWt; sw Ca | 72.11 |
| | | 1B | Soil TC, TN, S, TP; floc dryWt, AFDW; sw TP | 79.28 |
| | 2 | 2A | Soil TC, TP, AFDW; floc BD; sw TP | 77.99 |
| | | 2B | Soil TN, TP, BD, TC | 77.17 |
| | 7 | 5 | Soil AFDW, TC, TN, Ca; floc moisture | 77.90 |
| | 8 | 3 | Soil AFDW, Fe, macro AFDW, TC; floc AFDW | 77.87 |

342 * STA-1W Cell 1 PCA results consisted of similar and low average values, not
343 highlighting any particular driving variables.
344
345 **Table 5**: List of analyzed flow-ways by age, number of clusters and observed
346 spatial pattern of clusters (maps of cluster patterns available in supplementary
347 material).

| *STA* | *Flow-way* | *Oper. start year* | *Cell* | *No. clusters* | *Observed cluster pattern* |
|---|---|---|---|---|---|
| STA-1E | Central | 2006/7 | 3 | 4 | Zonal |
| | | | 4N | 3 | Zonal |
| | | | 4S | 3 | Zonal |
| STA-1W | Eastern | 1994 | 1 | 3 | Zonal gradient |
| | | | 3 | 4 | Zonal gradient |
| | Western | 1994 | 2 | 4 | Zonal gradient |
| | | | 4 | 4 | Tenuous zonal gradient |
| | Northern | 2000 | 5A | 4 | Zonal |
| | | | 5B | 4 | Tenuous zonal gradient |
| STA-2 | Flow-way 1 | 2000 | 1 | 5 | Zonal gradient |
| | Flow-way 2 | 2000 | 2 | 4 | Zonal |
| | Flow-way 3 | 2000 | 3 | 4 | Zonal gradient |
| STA-3/4 | Flow-way 1 | 2004 | 1A | 5 | Zonal |
| | | 2004 | 1B | 4 | Zonal |

| | | | | | |
|---|---|---|---|---|---|
| | Flow-way 2 | 2004 | 2A | 4 | Zonal |
| | | 2004 | 2B | 3 | Zonal |
| | Flow-way 3 | 2004 | 3A | 3 | Tenuous zonal gradient |
| | | 2004 | 3B | 5 | Zonal |
| STA-5/6 | Flow-way 1 | 1999 | 1A | 6 | Tenuous zonal gradient |
| | | 1999 | 1B | 3 | Tenuous zonal gradient |
| | Flow-way 2 | 1999 | 2A | 4 | Zonal |
| | | 1999 | 2B | 3 | Zonal |
| | Flow-way 7 | 1998 | 5 | 3 | Zonal gradient |
| | Flow-way 8 | 1998 | 3 | 5 | Zonal gradient |

348

349

350   CART analysis consistently found the strongest predictor variables for surface

351   water and floc TP to be other variables relating to P content (i.e. P in different

352   forms such as SRP, etc.) in soil, floc, and surface water; soil and floc BD; and soil

353   and floc TN. Measures relating to AFDW, TC and Ca also showed occasional

354   influence but were less widespread. Maps of CART model standard error by

355   location (not pictured) did not generally reveal any spatial relationships with

356   direction of flow, but did in some cases reveal zonal structures similar to the

357   cluster analysis (described below in 3.3).

358

359   Analysis with BBNs identified the strongest consistent predictors of recent year

360   NDVI to be variables relating to: vegetation type and cover, NDVI from previous

361   years, surface water TP, soil and floc TN, and soil and floc TC. BBNs predicting

362   surface water TP were most influenced by: other forms of surface water P, soil

363   BD, soil TN, soil TC, and soil TP.

364

365   **3.3. Spatial Trends**

366

367    Spatial patterns varied to a degree among treatment flow-ways. For instance, floc

368    and macrophyte characteristics dominated the models which predicted surface

369    water TP in STA-5/6; soil physical properties (e.g. bulk density) described many

370    of the spatial patterns in the treatment flow-ways of STA-3/4, etc.

371    Notwithstanding this, some general observations can be made regarding all

372    treatment flow-ways: (1) there are clear zonal patterns consistently present in

373    these systems that are, in many cases, independent of the direction of flow and

374    do not exhibit a simple linear gradient (Figure 2 shows STA-3/4 Flow-way 3 as

375    an example of purely zonal pattern; other examples include Flow-ways 1 and 2 in

376    the same STA and STA-1E's Central Flow-way, shown in supplementary

377    material); however these zonal patterns appear to align along the direction of

378    flow in the case of some older STAs and flow-ways (Figure 3 shows STA-1W's

379    Eastern flow-way as an example of zone-based gradient pattern; other examples

380    include STA-1W's Western Flow-way, STA-2's Flow-ways 1 and 3, and STA-5/6's

381    Flow-ways 7 and 8, shown in supplementary material and summarized in Table

382    5); (2) There is some consistency in the spatial arrangement of these zones over

383    the treatment flow-ways, such as surface water TP concentration being highest

384    close to the inflow structures and there closely associated with a zone of higher

385    floc and soil TP concentrations. Following these points, there is rarely any

386    further consistency in the spatial organization of zones, or in their

387    characterization, across flow-ways; but 3) soil TN often becomes an important

388    factor characterizing the zone around the outflow (e.g. STA-1W, STA-3/4).

389

390    **4. Discussion and Conclusions**

391    **4.1. Summary Statistics**

392

Two results stood out from the cell-wide summary statistics that were consistent with expectations. Firstly, the lowest mean values of internal surface water TP were found in flow-ways present in STAs 2 and 3/4, which have been previously cited as being two of the best-performing STAs for P removal (Pietro and Ivanoff, 2015). Secondly, all flow-ways consisting of multiple cells exhibited a trend of decreasing TP along the length of the flow-way (cell-wide summary statistics did not consider internal spatial patterns of single-cell flow-ways; these are discussed below), demonstrating the effects of P removal by the system at the STA scale. Taken broadly, this is consistent with the expectation that wetlands experiencing a uniform sheet flow should exhibit P decreases along a longitudinal flow-based gradient (Walker and Kadlec, 2011).

**4.2. Multivariate Analysis**

In considering the outputs from the data-mining analysis for the flow-ways; PCA is a general dimension reduction technique in which the underlying variation is maintained. It was used here because it is one of the primary steps in any multivariate data analysis as well as an effective way to represent variation in the data. Generally the PCA was successful, with an average of 75% of the variation explained. The most common variables identified as influential in the PC loadings were soil TC, soil TN, soil and floc BD, soil and floc TP and soil and floc AFDW. It should be noted that this particular analysis does not take into account non-continuous data (e.g. categorical variables such as soil series and parent material). In essence, the outcome from this analysis is an effective

417    summarization of the data but with little further insight into drivers, mainly

418    highlighting that most of the within cell/within flow-way variation is driven by

419    sediment nutrient concentrations and, to a lesser degree, floc TC and nutrient

420    content.

421

422    Cluster analysis resulted in cluster memberships that could be assigned to the

423    original data, revealing spatial patterns and structure in the data. Of interest here

424    were two points; do the data resolve clearly in clusters, and if so, how many (i.e.

425    how many classes of data are there in an STA flow-way), and are these classes

426    meaningful in any way? In general, most cells could be described by 3 to 5

427    clusters and only in one case (STA-5/6 Cell 1A; 6 clusters) were more clusters

428    needed (see Table 5). Clusters consistently grouped spatially into zone features

429    which did not appear to be tied to cell location within the flow path in many

430    cases; however in some cells these zonal features were observed to align along

431    the direction of flow. While not an unequivocal relationship, these 'zone-based

432    gradient' patterns appeared more likely to occur in older STAs and flow-ways

433    (Table 5). Patterns seemed only tenuously related to flow path at best in STAs-1E

434    and -3/4 (completed in 2007 and 2004, respectively), and generally more

435    obviously following the flow gradient in STA-1W (completed in 1994-2000),

436    STA-2 (completed in 2000), and STA-5/6 (completed in 1998/9).

437

438    The CART and BBN analyses both revealed similar relationships and driving

439    variables in the data. Surface water TP was found to share consistently strong

440    linkages with other forms of phosphorus in surface water (e.g. SRP and TSP) as

441    well as in floc and soil. Nitrogen, carbon, and bulk density in soil and floc also

442 factored in frequently; this highlights the potential importance of soil properties

443 to P dynamics in the STAs, as well as the possibility of coupled cycles wherein P,

444 N, and possibly C dynamics share co-dependencies and interrelationships.

445

446 **4.3. Observed Relationships and Drivers of P Dynamics**

447

448 It is evident from studies in the Everglades and elsewhere (Bayley and Mewhort,

449 2004; Bostic and White, 2007; Gu and Dreschel, 2008; Riggsbee et al., 2012), that

450 plant communities actively regulate P dynamics in wetlands. In the STAs, low

451 levels of water column P are achieved using strategic combinations of SAV and

452 EAV to address P in different forms and in different stages of the flow-ways

453 (Chen et al., 2015). In projecting this fact on the data mining exercise, one would

454 expect the spatial patterns of soil P to reflect plant community composition, and

455 plant communities would be expected to be a strong determinant in any

456 predictive model for soil or floc P. In our analysis this was only rarely the case;

457 however these effects may be obscured by the fact that much of the available

458 data on vegetation composition were categorical (e.g. vegetation class and

459 habitat type; NDVI being the notable exception as a continuous variable), and

460 thereby only possible to include in CART and BBN analyses. Both CARTs and

461 BBNs modeling surface water TP did not commonly reveal vegetation-related

462 measures as key predictors, but BBNs predicting NDVI frequently did highlight

463 surface water TP as an important driver (i.e. TP did not appear driven by

464 vegetation, but vegetation appeared driven by TP). Linkages between TP and

465 vegetation therefore may not be direct or omnipresent, but our analysis shows

466 support for some relationships.

467

468    Where P is limiting, or effectually buried, and therefore not available for the

469    plant communities, this may be reflected as plant stress (i.e. P limitation), which

470    can be remotely determined using NDVI (Henrik, 2012). The hypothesis is that

471    the indication of effective functioning of an STA is that, in the lower reaches of a

472    flow path, the vegetation may become P-limited. As a first instance, predictive

473    modeling of NDVI should indicate whether this is responsive to floc and soil

474    nutrient status. For BBNs predicting NDVI this indeed was the case; the strongest

475    predictors consistently included floc and soil nutrients, along with surface water

476    TP and other measures of vegetation health and composition. Note however that

477    prolonged exposure to low P concentrations may trigger a shift in plant

478    community composition to species that are more adapted to the low levels; such

479    a shift would be reflected in categorical habitat variables but not necessarily by a

480    decrease in NDVI. This highlights the importance of vegetation-related measures

481    beyond NDVI, and in turn the importance of methods such as BBNs that can

482    consider categorical expressions of vegetation community.

483

484    **4.4. Spatial Patterns of P and their Implications**

485

486    The observation that consistent spatial patterns appear zonal rather than based

487    on simple gradients is probably the most significant finding of the data mining, in

488    that the processes controlling P in these systems operate in zones in the

489    treatment flow-way, rather than along a smooth linear gradient as would be the

490    expectation (see Table 5). These zones are observed repeatedly across STAs and

491    flow-ways, and are consistently present as modeling outcomes (e.g. cluster

492   analysis and CART outputs) and as such are unlikely to be a modeling artifact.

493   There are a number of implications from approaching the STA flow-ways as

494   zones rather than a simple gradient. From a research perspective, the relative

495   importance of different factors, transformation and transport pathways of P

496   occurs in spatial patterns, and the form and shape of these patterns indicates the

497   relative importance of particular pathways. Likewise, this affects the

498   experimental sampling design, as these would then target zones rather than

499   seeking to measure along a gradient (biased sampling). From a management

500   perspective, this could simplify management options in that the operation and

501   management strategies can be directed at particular zones within a treatment

502   flow-way rather than an entire cell or the full flow-way, particularly once the

503   drivers of these zones are better understood. Nevertheless, in older STAs (e.g.

504   STA-1W, -2, and -5/6) these zonal patterns appeared to align more frequently

505   and obviously with the direction of flow, suggesting that P dynamics may

506   function largely in zonal patterns but slowly shift toward a zone-based gradient

507   pattern over the operational time of an STA. Of particular note, STA-2 flow-way 3

508   exhibited a strong gradient pattern in the cluster analysis result and has been

509   previously studied as one of the longest-running and best-performing treatment

510   flow-ways (Juston and Debusk, 2011; Juston et al., 2013).

511

512   The finding of zonal patterns of P concentrations in the STAs (whether forming

513   zone-based flow gradients or not), rather than simple uniform gradients

514   decreasing along the axis of water flow, differs from previous findings and the

515   usual expectation of P dynamics in wetlands (e.g. Kadlec, 1999; Walker and

516   Kadlec, 2011). One possible explanation for this difference is that the treatment

517  cells may be wide enough to allow partial mixing of water rather than a relatively

518  uniform sheet flow; this would account for more complex patterns (Walker and

519  Kadlec, 2011). If true, this would have implications for the assumptions made in

520  future flow modeling efforts in the STAs, and require a more complex

521  interpretation of the system than a one-dimensional sheet flow. Chen et al.

522  (2015) cautioned that analyses focused solely on inflow and outflow P

523  concentrations, while useful, do not consider P removal processes internal to the

524  treatment cells, as well as recommending that future studies consider

525  multivariate relationships. Doing so here has enabled additional findings, such as

526  the potential importance of relationships between P and soil factors, and the

527  possibility of P-N coupled cycles impacting dynamics. This latter result, while not

528  widely explored previously, is consistent with previous findings in Water

529  Conservation Area 2A (WCA 2a) on P and N functional linkages (White and

530  Reddy, 2003). Corstanje et al. (2009, 2007) found evidence that areas enriched

531  with P in WC-2a are mediated by N related parameters, such as potentially

532  mineralizable N and related microbial extracellular enzymatic activities. In STA

533  areas closest to the inflow, as P is relatively plentiful, the cycling P is likely to be

534  co-mediated by N and its dynamics.

535

536  **4.5. Data-Mining Advantages and Future Research**

537

538  Previous studies have examined the extensive data now available for P dynamics

539  in the STAs (e.g. Chen et al., 2015; Juston et al., 2013; Pietro and Ivanoff, 2015),

540  but this is one of the first known studies to comprehensively make use of the

541  diverse data collected in the interior treatment cells and flow-ways (e.g., soils,

542  vegetation, internal water quality) and the first to do so at such a broad scale

543  through a data mining approach. Doing so has facilitated new findings and

544  understanding around the functional P dynamics of the STA systems. Approaches

545  making use of these techniques are valuable for identifying biogeochemical

546  relationships, and should be considered and further employed in future studies

547  of the STAs as well as other engineered wetlands where sufficient data are

548  available.

549

550  In addition, there remain a number of further considerations moving forward.

551  First, many links between plant community composition and P dynamics remain

552  unclear beyond known differences between EAV and SAV in P removal (e.g.

553  Dierberg et al., 2002; Juston and DeBusk, 2006). In particular, we suspect there is

554  an element of scale effect; where these processes occur and are important at

555  scales finer than we considered in this study. Second, the approach used here

556  focused on data mining techniques, and while effective for exploring patterns in

557  the data it lacks a detailed process understanding of P biogeochemistry. The

558  incorporation of process understanding and process models (e.g. first order

559  equations) into the more stochastic modeling environment considered in this

560  study could produce a set of hybrid models which would both reflect process

561  knowledge and understanding but also, critically, allow for scaling and mapping.

562  Such an approach could better explore the process-based reasons for the zonal

563  patterns observed here and their potential relationships with flow-way age.

564  Finally, future research should seek to effectively consider the interaction

565  between different datasets available from the STAs in order to rigorously

566  consider time series analysis and pulsed events. A future study which initiates

567 with a thorough decomposition of the STA inflow and outflow data (volume and

568 concentrations), considers the stochasticity of this data and then moves to

569 incorporate it in the models of flow-way behavior should generate significant

570 insights in the STA dynamics, and to what degree performance is related to

571 stochastic events (e.g. storms or droughts) vs. deterministic processes (e.g. P

572 biogeochemistry, SAV, periphyton). Eventually this will relate to a measure of the

573 resilience of these systems; expressed as their capacity to withstand pressures

574 and maintain long term performance.

575

576 **4.6. Conclusions**

577

578 In conclusion, the use of data mining approaches on STA treatment cell and flow-

579 way data has identified, in a very general sense, spatial patterns in these systems.

580 These patterns are consistently zone-based across all flow-ways, which suggests

581 that the flow-ways function first as zonal systems rather than simple linear

582 gradient systems. Our analysis suggests that the primary drivers of the spatial

583 distribution of P in many of these systems are related to soil characteristics, and

584 that the zonal patterns of P distribution may begin to follow the predominant

585 flow path over time. The data further suggest the importance of coupled cycles in

586 these systems; in other words, the movement and transformation of P is coupled

587 to that of N.

588

589 **Acknowledgements**

592    **References**
593
594
595    Bayley, S.E., Mewhort, R.L., 2004. Plant community structure and functional

596        differences between marshes and fens in the southern boreal region of

597        Alberta, Canada. Wetlands 24, 277–294. doi:10.1672/0277-5212(2004)024

598    Bostic, E.M., White, J.R., 2007. Soil Phosphorus and Vegetation Influence on

599        Wetland Phosphorus Release after Simulated Drought. Soil Sci. Soc. Am. J.

600        71, 238. doi:10.2136/sssaj2006.0137

601    Breiman, L., Friedman, J.H., Olshen, R.A., Stone, C.I., 1984. Classification and

602        regression trees. Wadsworth, Belmont, California.

603    Chen, H., Ivanoff, D., Pietro, K., 2015. Long-term phosphorus removal in the

604        Everglades stormwater treatment areas of South Florida in the United

605        States. Ecol. Eng. 79, 158–168. doi:10.1016/j.ecoleng.2014.12.012

606    Corstanje, R., Portier, K.M., Reddy, K.R., 2009. Discriminant analysis of

607        biogeochemical indicators of nutrient enrichment in a Florida wetland. Eur.

608        J. Soil Sci. 60, 974–981. doi:10.1111/j.1365-2389.2009.01186.x

609    Corstanje, R., Reddy, K.R., Prenger, J.P., Newman, S., Ogram, A. V., 2007. Soil

610        microbial eco-physiological response to nutrient enrichment in a sub-

611        tropical wetland. Ecol. Indic. 7, 277–289. doi:10.1016/j.ecolind.2006.02.002

612    DeBusk, W.F., Reddy, K.R., Wang, Y., Koch, M.S., 1994. Spatial Distribution of Soil

613        Nutrients in a Northern Everglades Marsh: Water Conservation Area 2A. Soil

614        Sci. Soc. Am. J. 58, 543. doi:10.2136/sssaj1994.03615995005800020042x

615    Dierberg, F.E., DeBusk, T.A., Jackson, S.D., Chimney, M.J., Pietro, K., 2002.

616        Submerged aquatic vegetation-based treatment wetlands for removing

617        phosphorus from agricultural runoff: Response to hydraulic and nutrient

618        loading. Water Res. 36, 1409–1422. doi:10.1016/S0043-1354(01)00354-2

619    Diggle, P.J., Tawn, J.A., Moyeed, R.A., 1998. Model-based Geostatistics. Appl. Stat.

620        47, 299–350. doi:10.1111/1467-9876.00113

621    Entry, J.A., 2014. The Impact of Stormwater Treatment and Best Management

622        Practices on Nutrient Concentration in the Florida Everglades. Water, Air,

623        Soil Pollut. 225, 1758. doi:10.1007/s11270-013-1758-z

624    Gu, B., Dreschel, T., 2008. Effects of plant community and phosphorus loading

625        rate on constructed wetland performance in Florida, USA. Wetlands 28, 81–

626        91. doi:10.1672/07-24.1

627    Henrik, J.J., 2012. Utilizing NDVI and remote sensing data to identify spatial

628        variability in plant stress as influenced by management. Iowa State

629        University.

630    Ivanoff, D.B., Pietro, K., Chen, H., Gerry, L., 2013. Chapter 5: Performance and

631        Optimization of the Everglades Stormwater Treatment Areas, in: 2013 South

632        Florida Environmental Report - Volume I. South Florida Water Management

633        District, West Palm Beach, FL.

634    Juston, J., DeBusk, T.A., 2006. Phosphorus mass load and outflow concentration

635        relationships in stormwater treatment areas for Everglades restoration.

636        Ecol. Eng. 26, 206–223. doi:10.1016/j.ecoleng.2005.09.011

637    Juston, J.M., Debusk, T.A., 2011. Evidence and implications of the background

638        phosphorus concentration of submerged aquatic vegetation wetlands in

639        Stormwater Treatment Areas for Everglades restoration. Water Resour. Res.

640        47, 1–13. doi:10.1029/2010WR009294

641    Juston, J.M., DeBusk, T.A., Grace, K.A., Jackson, S.D., 2013. A model of phosphorus

642        cycling to explore the role of biomass turnover in submerged aquatic

643        vegetation wetlands for Everglades restoration. Ecol. Modell. 251, 135–149.

644        doi:10.1016/j.ecolmodel.2012.12.001

645    Kadlec, R.H., 1999. The limits of phosphorus removal in wetlands. Wetl. Ecol.

646        Manag. 7, 165–175. doi:10.1023/A:1008415401082

647    Kadlec, R.H., Wallace, S.D., 2009. Treatment Wetlands, 2nd ed. CRC Press, Boca

648        Raton, FL.

649    Malecki-Brown, L.M., White, J.R., Reddy, K.R., 2007. Soil Biogeochemical

650        Characteristics Influenced by Alum Application in a Municipal Wastewater

651        Treatment Wetland. J. Environ. Qual. 36, 1904. doi:10.2134/jeq2007.0159

652    McCormick, P. V., Rawlik, P.S., Lurding, K., Smith, E.P., Sklar, F.H., 1996.

653        Periphyton-Water Quality Relationships along a Nutrient Gradient in the

654        Northern Florida Everglades. J. North Am. Benthol. Soc. 15, 433–449.

655        doi:10.2307/1467797

656    McCune, B., Grace, J.B., 2002. Analysis of Ecological Communities. MjM Software,

657        Gleneden Beach, Oregon, USA.

658    Microsoft, 2003. Excel.

659    Newman, S., Kumpf, H., Laing, J.A., Kennedy, W.C., 2001. Decomposition

660        responses to phosphorus enrichment in an Everglades (USA) slough.

661        Biogeochemistry 54, 229–250. doi:10.1023/A:1010659016876

662    Norsys, 2014. Netica.

663    Orem, W., Newman, S., Osborne, T.Z., Reddy, K.R., 2014. Projecting Changes in

664        Everglades Soil Biogeochemistry for Carbon and Other Key Elements, to

665        Possible 2060 Climate and Hydrologic Scenarios. Environ. Manage. 55, 776–

666        798. doi:10.1007/s00267-014-0381-0

667    Pietro, K.C., Ivanoff, D., 2015. Comparison of long-term phosphorus removal

668        performance of two large-scale constructed wetlands in South Florida, U.S.A.

669        Ecol. Eng. 79, 143–157. doi:10.1016/j.ecoleng.2014.12.013

670    Reddy, K.R., DeLaune, R.D., 2008. Biogeochemistry of Wetlands: Science and

671        Applications. CRC Press, Boca Raton, FL.

672    Reddy, K.R., Jawitz, J.W., Paudel, R., Bhomia, R., Jerauld, M.A., 2009.

673        Comprehensive Analysis and Evaluation of Historical Data and Information

674        for the Stormwater Treatment Areas ( STAs ). Gainesville.

675    Reddy, K.R., Kadlec, R.H., Flaig, E., Gale, P.M., 1999. Phosphorus Retention in

676        Streams and Wetlands: A Review. Crit. Rev. Environ. Sci. Technol. 29, 83–

677        146. doi:10.1080/10643389991259182

678    Reddy, K.R., Newman, S., Osborne, T.Z., White, J.R., Fitz, H.C., 2011. Phosphorous

679        Cycling in the Greater Everglades Ecosystem: Legacy Phosphorous

680        Implications for Management and Restoration. Crit. Rev. Environ. Sci.

681        Technol. 41, 149–186. doi:10.1080/10643389.2010.530932

682    Richardson, C.J., 1999. The role of wetlands in storage, release, and cycling of

683        phosphorus on the landscape: a 25-year retrospective, in: Reddy, K.R.,

684        O'Connor, G.A., Schelske, C.L. (Eds.), Phosphorus Biogeochemistry of Sub-

685        Tropical Ecosystems. CRC Press, Taylor & Francis Group, Boca Raton, USA.

686    Riggsbee, J.A., Wetzel, R., Doyle, M.W., 2012. Physical and plant community

687        controls on nitrogen and phosphorus leaching from impounded riverine

688        wetlands following dam removal. River Res. Appl. 28, 1439–1450.

689        doi:10.1002/rra.1536

690    SAS, 2013. JMP.

691    South Florida Water Management District, 2015. South Florida Environmental

692        Report 2015. West Palm Beach.

693    South Florida Water Management District, Andreotta, H., Chimney, M., DeBusk,

694        T., Garrett, B., Gerry, L., Henry, J., Ivanoff, D., Jerauld, M., Kharbanda, M.,

695        Kirkland, M., Larson, N., Miao, S., Piccone, T., Pietro, K., Schwartz, L., Sierer-

696        Finn, D., Toth, L., Xue, S.K., Yan, Y., Zamorano, M., Zhao., H., 2015. Chapter

697        5B : Performance of the Everglades Stormwater Treatment Areas, 2015

698        South Florida Environmental Report. West Palm Beach.

699    StatSoft, 2014. Statistica 12.

700    Taalab, K., Corstanje, R., Zawadzka, J., Mayr, T., Whelan, M.J., Hannam, J.A.,

701        Creamer, R., 2015. On the application of Bayesian Networks in Digital Soil

702        Mapping. Geoderma 259-260, 134–148.

703        doi:10.1016/j.geoderma.2015.05.014

704    Walker, W.W., Kadlec, R.H., 2011. Modeling Phosphorus Dynamics in Everglades

705        Wetlands and Stormwater Treatment Areas. Crit. Rev. Environ. Sci. Technol.

706        41, 430–446. doi:10.1080/10643389.2010.531225

707    White, J.R., Reddy, K.R., 2003. Nitrification and Denitrification Rates of

708        Everglades Wetland Soils along a Phosphorus-Impacted Gradient. J. Environ.

709        Qual. 32, 2436. doi:10.2134/jeq2003.2436

710

711

712     **Figures**



Figure 1: Locations of the Stormwater Treatment Areas in south Florida, USA, indicating individual treatment cells and direction of flow. Bolded flow-way names and darkened arrows denote flow-ways included in analysis.

Figure 2: Spatial patterns detected by cluster (A) and CART (B) analyses – an

example for STA-3/4 flow-way 3. Image B represents the distribution of

CART nodes (symbol numbers represent the number of nodes in the CART

model) corresponding to the prediction of surface water total P

(concentration denoted by symbol color). Note that patterns are

predominantly zonal and only tenuously aligned with flow direction.

Figure 3: Spatial patterns detected by cluster (A) and CART (B) analyses – an example for STA-1W Eastern flow-way. Image B represents the distribution of CART nodes (symbol numbers represent the number of nodes in the CART model) corresponding to the prediction of surface water total P (concentration denoted by symbol color). Note that zonal patterns appear largely aligned with flow direction, indicating a gradient-based behavior to the individual zones.

**A datamining approach to identifying spatial patterns of P forms in the Stormwater Treatment Areas in the Everglades, US**

Corstanje, R., Grafius, D.R., Zawadzka, J., Moreira J., Vince, G., Ivanoff, D., Pietro, K.

**Supplementary Materials**

These maps show the K-mean clusters and tree nodes resulting from Classification and Regression Trees (CARTs) analysis performed within particular flow-ways of the Stormwater Treatment Areas (STAs) that had sufficient data to do so. Please note that, in the case of CARTs, the results are only shown for the flow-ways with availability of data on total surface water phosphorus.

**Figure S1**: Cluster analysis for STA-1E Central flow-way. Arrows indicate the direction of water flow through the flow-way.

**Figure S2**: Results of A – cluster analysis, and B – CARTs analysis for STA-1W Eastern flow-way. Arrows indicate the direction of water flow through the flow-way. Numbers in CART results indicate the number of nodes in the CART model.
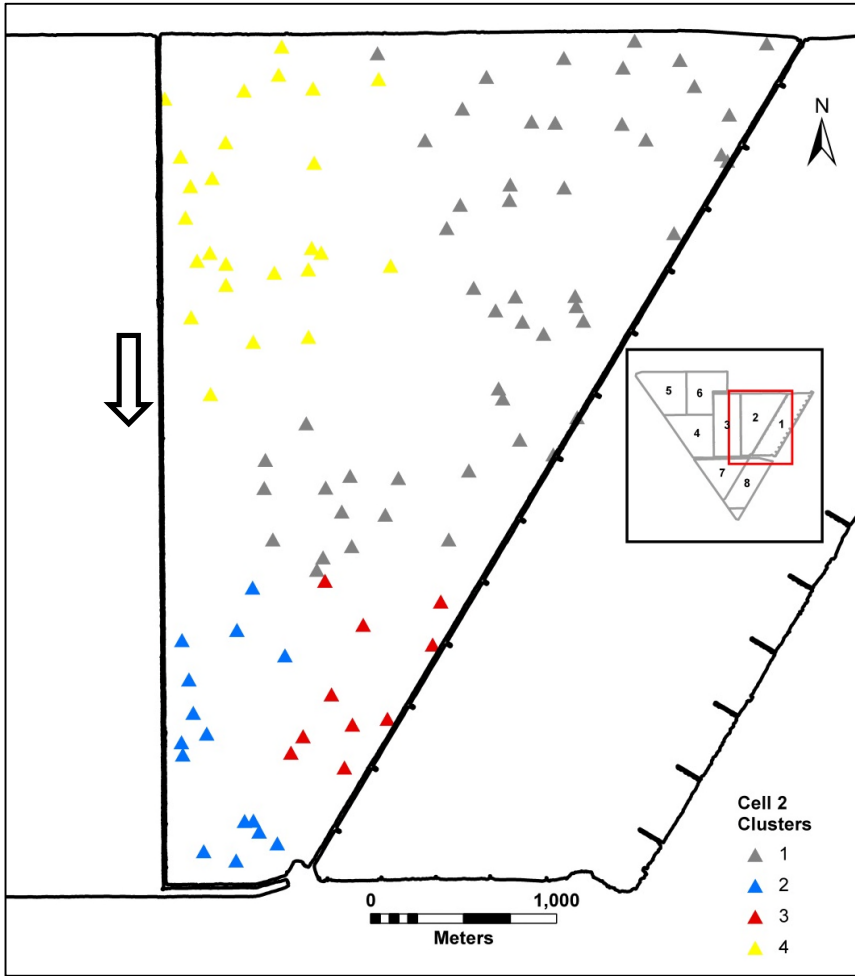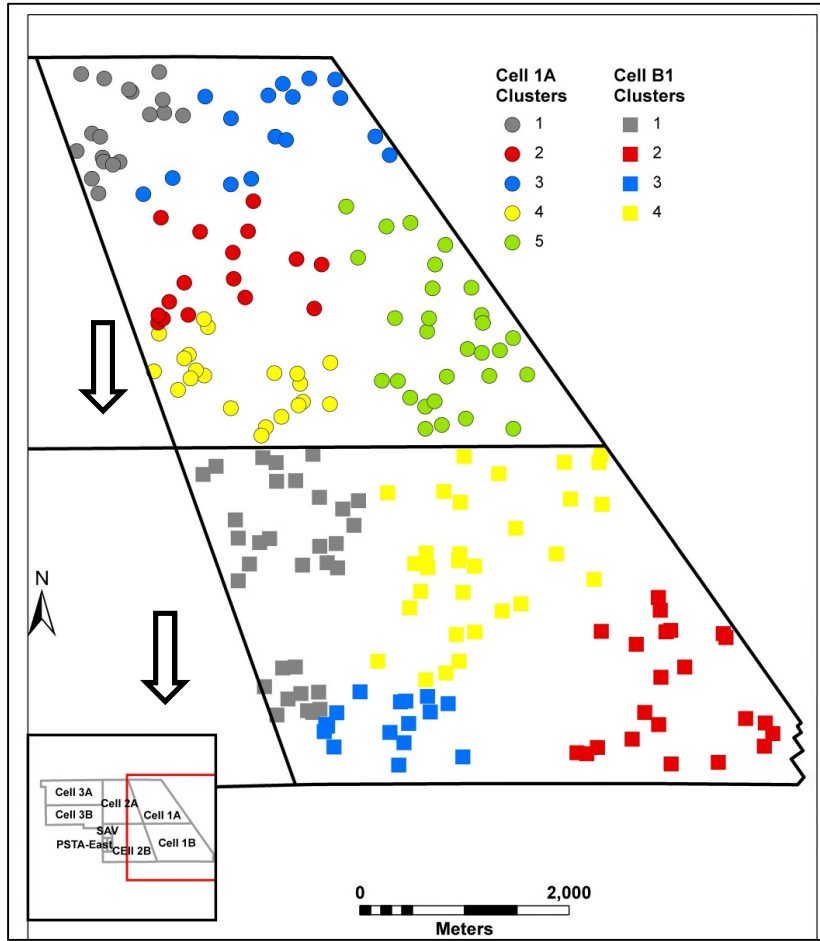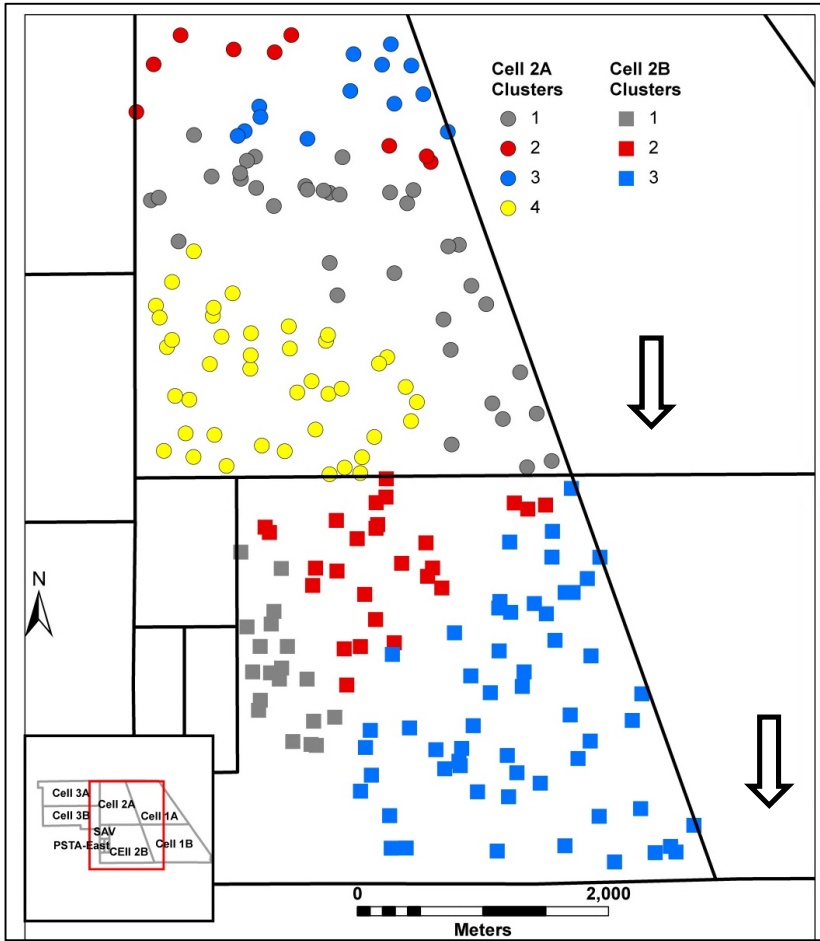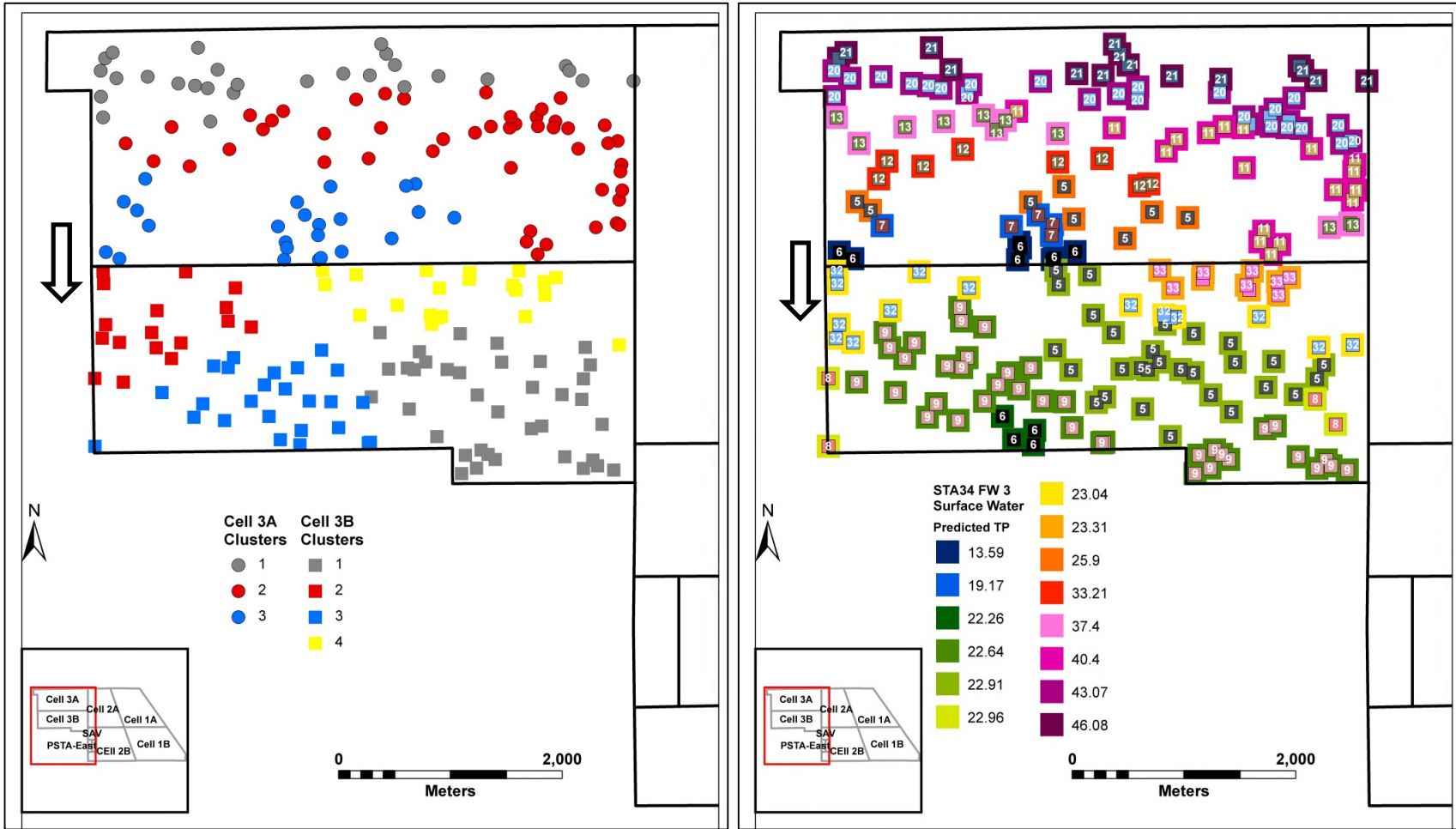
**A**  **B**

**Figure S2:** Results of A – cluster analysis, and B – CARTs analysis for STA-1W Western flow-way. Arrows indicate the direction of water flow through the flow-way. Numbers in CART results indicate the number of nodes in the CART model.

**Figure S3:** Results of A – cluster analysis, and B – CARTs analysis for STA-1W Northern flow-way. Arrows indicate the direction of water flow through the flow-way. Numbers in CART results indicate the number of nodes in the CART model.

**Figure S4:** Results of A – cluster analysis, and B – CARTs analysis for STA-2 flow-way 1. Arrows indicate the direction of water flow through the flow-way. Numbers in CART results indicate the number of nodes in the CART model.
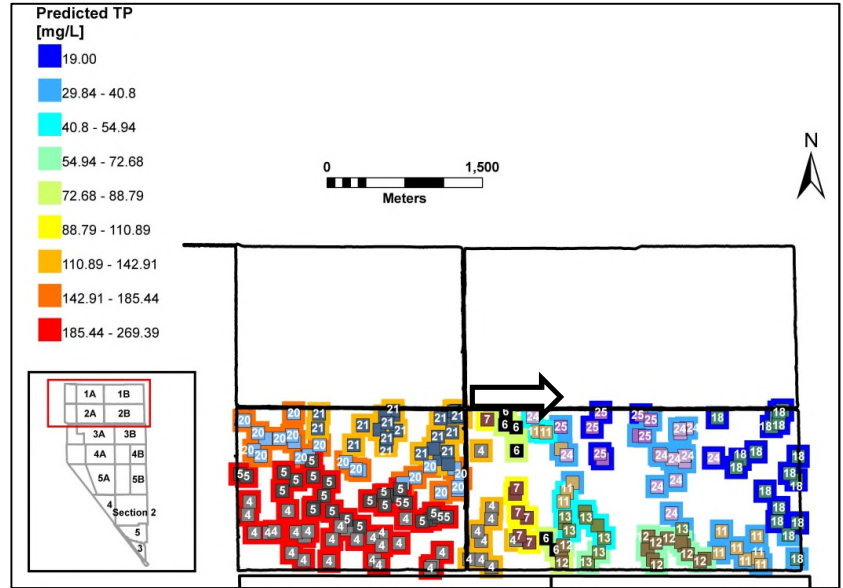
**A**

**B**

**Figure S5:** Results of A – cluster analysis, and B – CARTs analysis for STA-2 flow-way 3. Arrows indicate the direction of water flow through the flow-way. Numbers in CART results indicate the number of nodes in the CART model.

**Figure S6:** Cluster analysis for STA-2 flow-way 2. Arrows indicate the direction of water flow through the flow-way.

**Figure S7**: Results of cluster analysis for A – STA-3/4 flow-way 2 and B – STA-3/4 flow-way 1. Arrows indicate the direction of water flow through the flow-way.
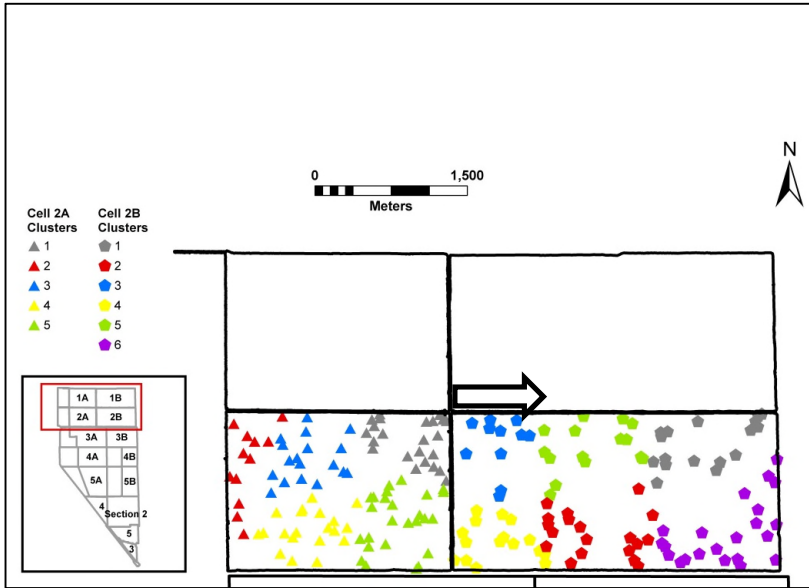
**Figure S8:** Results of A – cluster analysis, and B – CARTs analysis for STA-3/4 flow-way 3. Arrows indicate the direction of water flow through the flow-way. Numbers in CART results indicate the number of nodes in the CART model.
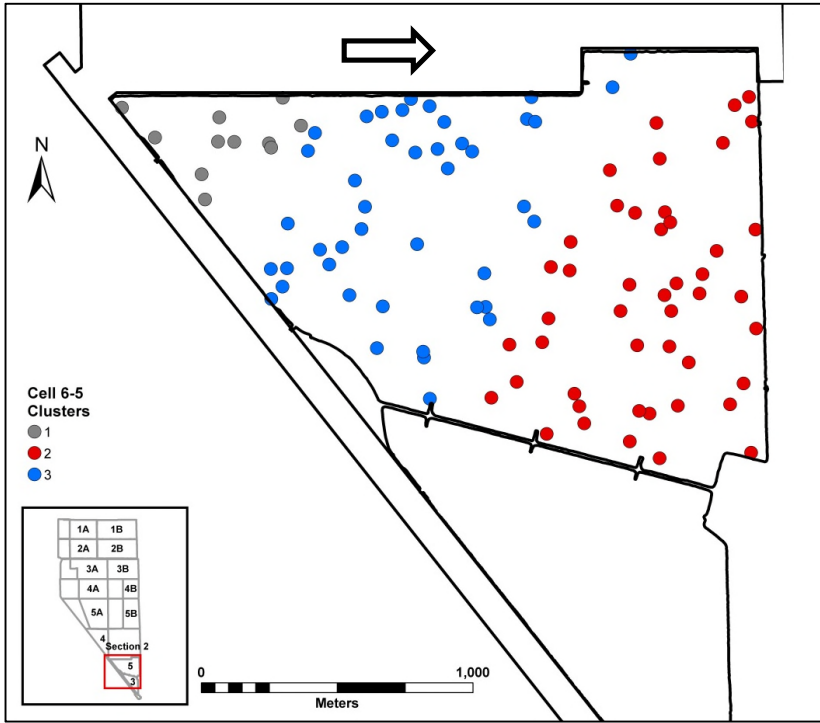
**A**

**B**

**Figure S9**: Results of A – cluster analysis, and B – CARTs analysis for STA-5/6 flow-way 1. Arrows indicate the direction of water flow through the flow-way. Numbers in CART results indicate the number of nodes in the CART model.
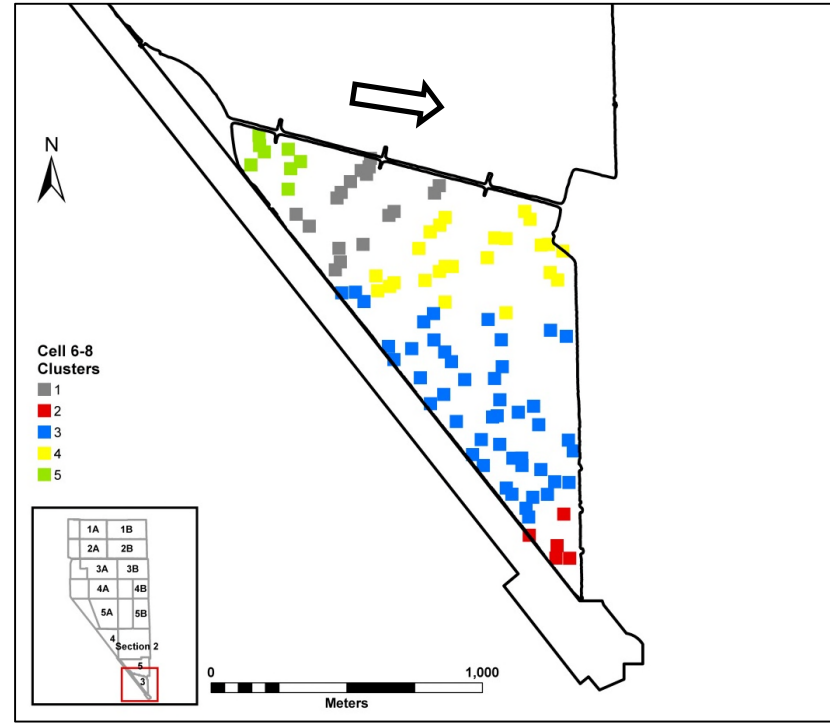
**A**

**B**

**Figure S11:** Results of A – cluster analysis, and B – CARTs analysis for STA-5/6 flow-way 2. Arrows indicate the direction of water flow through the flow-way. Numbers in CART results indicate the number of nodes in the CART model.

**A**

**B**

**Figure S10:** Results of cluster analysis for A – STA-5/6 flow-way 7 and B – STA-5/6 flow-way 8. Arrows indicate the direction of water flow through the flow-way.