

CVNetica—A cross-validation package driving Netica with Python

Michael N. Fienen

*US Geological Survey, Wisconsin Water Science Center, 8505 Research Way, Middleton WI
53562 USA*

Nathaniel G. Plant

*US Geological Survey, St. Petersburg Coastal and Marine Science Center, 600 Fourth Street
South, St. Petersburg, Florida, 33701, USA.*

Abstract

1 Bayesian networks (BNs) are powerful tools for probabilistically simulating nat-
2 ural systems and emulating process models. Cross validation is an important
3 technique to avoid overfitting that can result from overly complex BNs. Over-
4 fitting results in a reduction of true predictive skill. Formal cross-validation
5 procedures to evaluate relative performance in prediction outside of calibration
6 data has been discussed in various studies but rarely implemented. The lack of
7 widespread cross-validation is due in part to a lack of software tools designed to
8 work with available BN packages. CVNetica is an open-source package written
9 in Python that extends the Netica software package to perform cross-validation
10 and includes a framework to read, rebuild, and learn BNs from data. Optimal
11 BN complexity can be identified through exploration of numbers of bins, nodes,
12 and edges making up a BN using predictive skill as a metric of performance.
13 The insights gained from cross-validation and implications on predictive versus
14 descriptive skill are illustrated with two examples: a data-driven oceanographic
15 application in which wave height predictions are made from time series data
16 learned from nearby buoys and weather stations; and a model-emulation appli-
17 cation in which the source of water to groundwater wells is evaluated using a

18 BN trained to the results of a large, regional groundwater process model.

1. Introduction

19 Over the past two decades, the use of Bayesian Networks [BN; 1] has in-
20 creased greatly, in large measure due to the availability of commercial software
21 packages such as Netica [2] and Hugin [3] among many others. Applications
22 in water resources have included groundwater management [4, 5, 6], and model
23 emulation [7, 8, 9]. This builds on a history of applications in national security,
24 economics, and ecology.

25 An important topic that is not always discussed in the literature is that
26 applications of BNs need formal tests and validation of prediction performance
27 [10, 11]. Some validation metrics are calculable by the commercial software
28 packages, but substantial gaps in capabilities remain. Fortunately—at least in
29 the case of Netica—an application programming interface (API) exists with ver-
30 sions in multiple programming languages. To create a toolbox of performance
31 metrics, we used Python [12] with the Netica C APIs. These APIs expose most
32 of Netica’s functionality, through functions, to external programming. Among
33 the languages available, C was chosen because one of our goals was to interface
34 with Python 2.7.6 [12], Numpy 1.8 and Scipy 0.13.2 [13]. We discuss the tech-
35 nical challenges associated with running C APIs using Python and describe the
36 toolbox of validation metrics included in this work.

37 Building on techniques introduced by Fienen et al. [9], we developed tools
38 addressing two fundamental questions of Bayesian network performance: how
39 does predictive performance compare with descriptive calibration quality?; and
40 how does the complexity of the underlying network impact predictive and de-
41 scriptive performance? Cross-validation is used to answer both questions, and
42 the number of bins per node is used as a metric of complexity to answer the

second. These specific questions are evaluated in this work but our framework allows for consideration and analysis of other validation metrics and techniques beyond those presented here.

2. Bayesian Networks

This background section on Bayesian networks (BNs) is derived from Fienien et al. [9]. A Bayesian network is a directed acyclic graph [14], composed of nodes and edges. Nodes represent variables whose parameter values may include Boolean, discrete states, or, for continuous variables, discrete ranges that are discretized into bins. Edges form the connections between nodes and represent a correlated connection between the properties represented by the nodes. The entire catalog of these correlations make up conditional probability tables (CPTs). In a predictive context, nodes can be thought of as either input (e.g. forcing) or output (e.g. response), although this distinction is not a sharp one as the correlations learned by the BN are ambivalent with respect to direction. Nodes can also be intermediate if they act as constraints or model coefficients.

An example of a BN created and visualized using Netica is presented in figure 1.

Calculations are made using the BN based on conditional probabilities using Bayes' Theorem

$$p(F_i|O_j) = \frac{p(O_j|F_i)p(F_i)}{p(O_j)} \quad (1)$$

where $p(F_i|O_j)$ is the posterior (updated) probability of a forecast (F_i) given (conditional on) a set of observations (O_j); $p(O_j|F_i)$ is the likelihood function, $p(F_i)$ is the prior probability of the forecast, and $p(O_j)$ is a normalizing constant. The posterior probability reflects an updating that is achieved by considering the entire chain of conditional probabilities of all bins connected to the node representing F_i . The likelihood function represents the probability that the observations (O_j) would be observed given that the forecast was perfectly

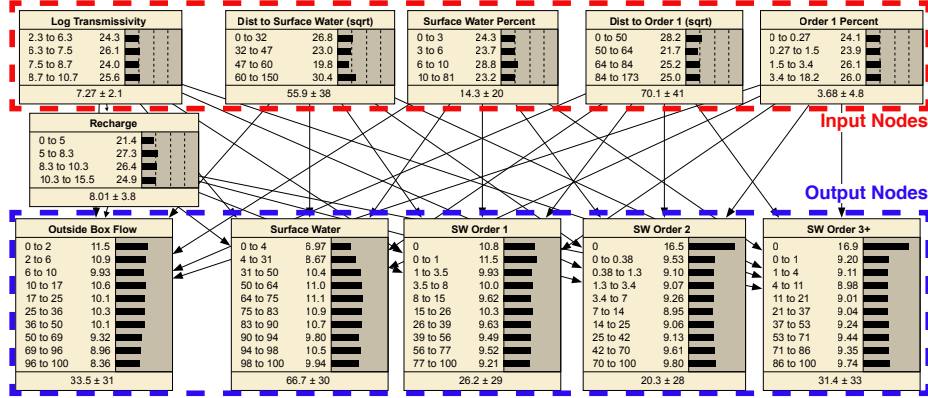


Figure 1: Example groundwater application of a Netica BN showing input (outlined in a red box) and output (outlined in a blue box) nodes, edges (black lines) and, in this case, a single intermediate node (recharge).

68 known. This is a metric of the ability of the BN to function as a forecasting de-
69 vice and imperfections in such forecasts are a function of epistemic uncertainty.
70 Epistemic uncertainty includes uncertainty due to model imperfection, data er-
71 rors, data paucity, and other sources. The prior probability of the forecast,
72 $p(F_i)$, is the probability of a forecast without the benefit of updated obser-
73 vations and the BN (or a process model or other experiment). $p(F_i)$ may be
74 calculated by using expert knowledge, or may be assumed relatively uninforma-
75 tive to make the entire process as objective as practical (similar to an ignorance
76 prior [15]). A common prior often used in BNs is the division of a node into
77 bins of equal probability. This results in bins of equal probability or “belief”
78 although it is not exactly an ignorance prior because the probability mass in
79 each bin may differ due to variable bin widths. It is possible to evaluate the
80 contribution to all uncertainty values calculated by the BN by expressing the
81 uncertainty in the prior probabilities. In Figure 1, the horizontal bars corre-
82 spond to relative probabilities associated with bins outlined by the numbers
83 listed to the left of them. These bars form a histogram and are referred to as

84 “belief bars.”

85 Once a system is cast in a BN, new observations of system state are applied
86 and propagated through the BN using Bayes’ theorem such that all forecasts
87 made in the model are contingent upon the specific observations of system state.
88 In other words, each forecast is associated with a specific configurations of sys-
89 tem state. In our approach, observations are indicated by selecting a bin and
90 forcing the probability of a value in the node to be 100%. This implies that
91 observational uncertainty does not exceed the width of the specified bin (for
92 continuous variables) or that the discrete or Boolean state is known perfectly.
93 (It is straightforward to relax this assumption to consider inputs that are un-
94 certain.) When this operation is performed, the Bayesian update propagates in
95 each direction among nodes that are d-connected [1], updating the probabilities
96 regardless of causal direction. In this way, correlations are expressed as well
97 as causal responses. By selecting a suite of observations of state, the BN acts
98 like a transfer function by providing an estimate of the forecast of interest and
99 associated uncertainty.

100 A key piece of a priori information is the establishment of edges connecting
101 the nodes. Edges should reflect a cascade of causality grounded in an under-
102 standing of the underlying process being modeled. If multiple processes from
103 different models are to be linked, the selection of edge relationships defines the
104 linkage. While machine learning can be used to teach a BN which parameters
105 are connected to each other and to outputs, we adopt a Bayesian approach in
106 which expert system understanding is used to specify these connections through
107 the identification of nodes and edges. In this way, the BN honors the physical
108 conditions known by the modeler and these are incorporated as soft knowledge.

109 In Figure 1, arrows on the edges indicate the direction of causal dependence.
110 When all nodes are d-connected, the direction of the edge arrows serve no pur-

pose. However, in the context of d-separation, the direction of causality has important ramifications on the propagation of uncertainty from observations to forecasts.

When computational conditions and problem size permit, a conditional probability table (CPT) can be created that directly enumerates the conditional probabilities of all nodes in the BN. This becomes impractical rapidly, however, because the size of the CPT scales on the order of $n \times d^{k+1}$ where n is the number of nodes, d is the number of bins, and k is the number of parents for a node. In the case where full enumeration is impractical due to this rapid increase in computational expense with complexity, an iterative expectation-maximization (EM) algorithm is used [16] to calculate approximate probabilities and maximum-likelihood values for the BN without full enumeration of the CPT. The EM algorithm iterates between estimating the maximum log likelihood of the function and finding the set of parameters resulting in that maximum log likelihood.

3. Cross validation tool

CVNetica is a Python module that performs cross-validation and calculates other performance metrics on BNs created with the Netica software package. Netica is a commercial package with more power than open-source alternatives. However, CVNetica is open-source and freely available. The APIs for Netica are described in [17] and are provided as a dynamic linked library (DLL) for Windows. Static libraries are also available for Macintosh and *nix platforms, but to use them with Python, dynamic interface wrappers would be necessary in addition to the Python function wrappers written in CVNetica.

The core functionality of CVNetica is based around the concept of using cross-validation [18, 11, 9] metrics to assess the quality of predictions made by a BN. In k -fold cross validation used in this work, the calibration dataset is, ran-

domly without replacement, divided into k folds or partitions where k typically is between 2 and 10. For each fold, the BN is trained using the dataset without the data in the fold, then the BN is used to make predictions on the left-out data. In this way, performance of the BN is evaluated on data not used in calibration to simulate performance in true future prediction. Several performance metrics can be used for this purpose, as discussed in Norsys Software Corp. [2], Plant and Holland [7], and Fienen et al. [9]. In this work, we will focus on skill

$$sk = \left[1 - \frac{\sigma_e^2}{\sigma_o^2} \right] \quad (2)$$

where σ_e^2 is the mean squared error between observations and BN predictions, and σ_o^2 is the variance of the observations [19, 7, 20]. Skill is evaluated by comparing BN predictions to observations with a value of unity indicating perfect correspondence and a value of zero indicating substantial discrepancy between BN predictions and observations.

CVNetica also reports log loss, error rate, experience, quadratic loss, mutual information (entropy), variance reduction (sensitivity) all of which are described by Norsys Software Corp. [2, 17]. Expected values are reported either as mean or most likely (ML). For ML values, the value corresponding to the center of the bin with the highest predicted probability is reported. The mean values, are computed as the product of the bin centers and the probability in each bin, consistent with a typical expectation operation.

By evaluating skill over both the calibration data sets and prediction data sets, the value of a BN as a descriptive or predictive tool can be evaluated. As BN complexity increases, so does the calibration sk and with sufficient complexity, calibration sk approaches unity (perfection). However, greater descriptive value in a BN comes at a cost in predictive value. This is the classic condition of overfitting as cast in the context of information theory by Fienen et al. [9].

162 One way to systematically evaluate BN complexity is to adjust the number
 163 of bins for each node with more bins meaning a greater level of complexity. CV-
 164 Netica has the capability to make this type of analysis efficient by allowing the
 165 user to specify an original BN and a configuration of bins for each node. CV-
 166 Netica then builds a new BN with the requested number of bins and assigning
 167 equiprobable prior distributions for each bin. In Fienen et al. [9] the number
 168 of bins was assumed the same for each node. Using CVNetica the number of
 169 bins in each node can be varied independently to allow for exploration of vari-
 170 ous assumptions of complexity. The user can also establish scenarios manually
 171 varying the number and nature of edges connecting nodes and even the number
 172 of nodes themselves. A group of these scenarios is defined by CVNetica as a
 173 “set.” Each set can be evaluated as a batch and then tabulated and graphical
 174 results are generated of performance metrics across the sets.

3.1. *Details about program structure*

175 There are two levels at which CVNetica performs. At the highest level, a
 176 script in `CV_driver.py` performs the cross validation protocol described below.
 177 This script is driven by an XML-based configuration file and should generally
 178 require minimal editing, save for identifying the configuration file to use in the
 179 `parfile` variable name. At a lower level, `pythonNeticaTools.py` provides the
 180 `pyneticaTools` class that interacts with the Netica DLL via wrappers around
 181 many essential Netica functions. Examples of how these methods work are
 182 discussed in Section 4.1. At an intermediate level, `pythonNetica.py` provides
 183 the `pynetica` class that combines several Netica functions for tasks such as
 184 starting a Netica environment, rebinning nodes, and other intermediate level
 185 tasks.

186 3.1.1. Cross Validation Driver

187 The `CV_driver.py` script drives a cross-validation exercise specified in the
188 XML based configuration file (Figure 2). If no rebinning is requested (`< rebin_flag > False < /rebin_flag`
189 the BN specified in the `baseNET` element is used for analysis along with the
190 casefile identified by the `baseCAS` element and metrics of performance. If the
191 `rebin_flag` element is True, then the nodes from the BN identified in the
192 `originalNET` element are rediscritized using the information on rebinning pro-
193 vided at the end of the input file. For each node listed, if `numbins > 0` the
194 node is discretized into bins `numbins` bins of equal probability. In the special
195 case where `numbins = 0`, the node is not rediscritized but it is used either as
196 input or response as described by the `input` and `response` elements above.
197 This special case allows for other discretization strategies (such as thresholds)
198 to be implemented for nodes that are to be treated as input or response nodes
199 but without equiprobable discretization. Nodes that are not identified as either
200 input or response should not have `node` elements provided and are unaltered by
201 CVNetica in the analysis.

202 If the `CVflag` element is False, only a single run using all the data in the
203 `baseCAS` file and the BN identified in the `baseNET` is performed and metrics are
204 calculated. The predictions for each configuration of input are recorded in a
205 compressed Python pickle file.

206 If the `CVflag` element is True, then k-fold cross validation is performed using
207 the number of folds indicated in the `numfolds` element. For each fold, $\frac{n}{k}$ (where
208 n is the total number of data points and k is the number of folds) data points
209 are separated from the rest of the data points to be left out of the calibration,
210 selecting from a randomized list such that each fold samples across the training
211 set to span spatial or temporal trends or patterns. The BN is then retrained on
212 the $n - \frac{n}{k}$ retained data and metrics of performance are calculated for both the

```

<data>
  <control_data>
    <baseNET>glacial_bins4_5_0.neta</baseNET> <!-- name of main .neta file -->
    <baseCAS>glacial.cas</baseCAS> <!-- name of main data file -->
    <rebin_flag>True</rebin_flag> <!-- flag determining if
                                rebinning should be performed -->
    <originalNET>glacial.neta</originalNET> <!-- original .neta file providing node
                                structure and bins of numbins=0 below-->
    <pwdfile>mikeppwd.txt</pwdfile> <!-- name of Netica license file -->
  </control_data>
  <kfold_data>
    <CVflag>True</CVflag> <!-- flag indicating if k-fold cross validation
                                should be carried out -->
    <numfolds>10</numfolds> <!-- number of folds for cross validation -->
  </kfold_data>
  <scenario>
    <name>glacial_set1</name> <!-- scenario name for output files -->
    <input>sqrT_SW_MIN</input> <!-- input tags identify nodes as used for input -->
    <input>sqrT_RIVMIN1</input>
    <input>PCTORD1</input>
    <response>EXT_FLOW</response> <!-- response tags identify nodes as used for output -->
    <response>SW_SRC</response>
  </scenario>
  <sensitivity>
    <report_sens>True</report_sens> <!-- flag indicating if Netica sensitivity and
                                other built-in metrics should be reported -->
  </sensitivity>
  <learnCPTdata>
    <voodooPar>100</voodooPar> <!-- fitting parameter for learning CPTs -->
    <useEM>True</useEM> <!-- use EM to learn CPTs if True. Else, use
                                incorporate casefile method -->
  </learnCPTdata>
  <rebinning>
    <!-- if rebin_flag is True, then bin_setup.py will read in the
    rebin_name to write out the rebinned .neta file and will
    use the newbins information for that purpose.
    Nodes will be rediscritized into numbins equiprobable bins.
    Special case when numbins = 0, the node is not rediscritized from originalNET -->
    <newbins>
      <node numbins="4">sqrT_SW_MIN</node>
      <node numbins="4">sqrT_RIVMIN1</node>
      <node numbins="4">PCTORD1</node>
      <node numbins="5">EXT_FLOW</node>
      <node numbins="0">SW_SRC</node>
    </newbins>
  </rebinning>
</data>

```

Figure 2: Example XML configuration file for defining problem parameters. Blue text identifies syntax of element names, green text indicates comments in the file, and bold black text indicates element values. In the special case of the `node` element, an attribute (`numbins`) is indicated in red text.

213 left out data (referred to as “validation”) and the training data (referred to as
 214 “calibration”).

4. Working with Ctypes

215 The Netica software provides APIs for accessing and using the functions
 216 within it. Several versions of these APIs are available as precompiled libraries.
 217 To interface with Python, the C programming language APIs can be interfaced
 218 using the `ctypes` module which is built-in to Python 2.5+. The `ctypes` mod-
 219 ule enables the use of functions from a dynamic library of C code (a DLL on
 220 Windows) in the Python environment. In addition to making the functions

221 accessible, some translation of variables is required—for example, C often refers
222 to data using pointers whereas Python does not explicitly do so. C functions
223 often return pointers to memory space of the resulting arrays so `ctypes` must
224 be used to read the correct amount of data from memory to populate an array
225 for further use in Python.

226 CVNetica provides Python functions wrapped around Netica C functions
227 and helper functions to translate data to and from the Python environment. In
228 the remainder of this section, the main aspects of interfacing with the Netica
229 APIs are discussed in general terms. These examples use code snippets from
230 the CVNetica codebase. Further documentation about `ctypes` is available from
231 the official documentation (<http://docs.python.org/2/library/ctypes.html>).

4.1. Accessing the DLL

232 The first task when accessing the Netica DLL is to make the functions avail-
233 able to Python by assigning the DLL to an object. Note that the filename is not
234 in quotes, nor is the `.dll` extension required. The `ctypes` module is imported
235 as `ct` so in future code descriptions, `ct.<>` implies a method or property from
236 `ctypes`.

```
237 import ctypes as ct
238 self.n = ct.windll.Netica
```

239 After this, `self.n` is an object with all of the Netica API functions avail-
240 able. To call a function from the DLL, the function name is dereferenced from
241 `self.n` and in CVNetica, a wrapper function is created as an interface to the
242 Netica function. In the following example, the Netica function to be called
243 is `EnterNodeValue_bn`. This function takes two arguments as indicated in
244 the function definition by Netica: `void EnterNodeValue_bn (node_bn* node,`
245 `double value)` [17]. The two arguments are of the custom C type defined by
246 Netica as `node_bn*` `node` and a double-precision float `double value`. A wrap-

per around this function must then make type conversions as appropriate. The CVNetica variable `cnode` was returned by a Netica function, so it is already of the type required (a pointer). However, the CVNetica variable `cval` is a Python float and must be converted to a C double using a `ctypes` conversion.

```
def EnterNodeValue(self, cnode, cval):
    self.n.EnterNodeValue_bn(cnode, ct.c_double(cval))
    self.chkerr()
```

The `chkerr` method polls the Netica DLL for current error status and, if an error is encountered, kills CVNetica and displays the error from Netica to standard error.

4.2. Exchanging information with the Netica DLL

The functions in Netica can accept a variety of argument types. In the `pyneticaTools` class, methods that function as wrappers around Netica functions are written. The names are the same as the Netica functions with the `_bn`, `_cs`, and `_ns` suffixes removed. This class is not specific to cross validation applications and is meant to also serve as a starting point for other applications in which Netica functions must be used in Python.

The easiest type is a pointer to an object returned by another Netica function. In this case, a Python variable represents the pointer—just a memory address—so no conversion is necessary. For single Python floats and ints, the conversions are `ct.c_double(cval)` and `ct.c_int(cval)`, respectively, where `cval` is the Python variable.

Some Netica functions return a double value but also write another result to memory at a location indicated by a pointer passed to the function. An example is `GetNodeExpectedValue_bn`. The structure of this function in C is

```
double GetNodeExpectedValue_bn (node_bn* node,
```

273 double* std_dev , double* x3 , double* x4)

274 where the returned value is the expected value (double precision) of the node
275 identified by `node_bn*`, the standard deviation is written to the memory location
276 identified by the pointer `double* std_dev`, and `x3` and `x4` are NULL pointers
277 reserved for future implementation. To collect the main returned value of the
278 function, we must set `restype` of the function—accomplished through making an
279 alias temporary function—and accepting the value as normally with a function.
280 To make use of the returned second value in Python—the value written to a
281 memory location identified by a pointer—we must pass a `double` variable by
282 reference (in other words, a pointer to the `double`). The Python wrapper for
283 `GetNodeExpectedValue_bn` illustrates this process

284

```
285 def GetNodeExpectedValue(self , cnode):  
286     std_dev = ct.c_double()  
287     tmpNeticaFun = self.n.GetNodeExpectedValue_bn  
288     tmpNeticaFun.restype=ct.c_double  
289     expected_val = tmpNeticaFun(cnode , ct.byref(std_dev) ,  
290                               None , None)  
291     self.chkerr()  
292     return expected_val , std_dev.value
```

293 Some Netica functions return either a character array or a numerical array.
294 In both cases, the C code in Netica returns a pointer to the data. The Python
295 code must, then, read a specified amount of data from that pointer location.
296 Unlike pure Python, it is possible to read off the end of the information starting
297 at the pointer location, so we must also specify the number of values to read
298 from the memory location. Helper functions in `cthelper.py` read the character
299 pointers, and single and double precision pointers. An example of this being

300 used in CVNetica is in the `ReadNodeInfo` method of the `pyneticaTools` class.

5. Example Applications

301 The CVNetica code was applied to two different applications to evaluate
302 predictive performance and guide the appropriate level of complexity for BN
303 design. The two applications are (1) a data-driven prediction of ocean wave
304 evolution and (2) a model emulation using a BN to make predictions trained on
305 results of a groundwater flow model.

5.1. Data driven ocean waves

306 Weather forecasting and modeling has achieved sufficiently high accuracy
307 that it is possible to replace observations with models if models are initialized
308 well and have good boundary condition data. However, weather forecasts are
309 not routinely available for periods extending more than a few days ahead, and
310 they become less accurate. We would like to allow the climatological prior in-
311 formation to inform predictions when observations or forecasts are not available
312 or are uncertain. As an example, we would like to predict wave height just
313 offshore of the coast where there are not persistent observations. This could be
314 done with laborious Monte Carlo simulations using models and previous clima-
315 tology for model initialization and boundary-condition forcing. Or, we could use
316 extant model output or observations to learn both the sensitivity of a specific
317 prediction to changes in boundary conditions and include uncertainty in this
318 sensitivity (the joint correlation) as well as uncertainty in the boundary con-
319 ditions. This approach has been implemented before using Bayesian networks
320 [7, 8], but the fidelity of the of the resulting BN models was not examined in de-
321 tail, other than to note that to maximize the consistency of the BN predictions
322 with new observations, the BN inputs should include input parameter-value or
323 data uncertainties. Were these BN models over-fit to the data?

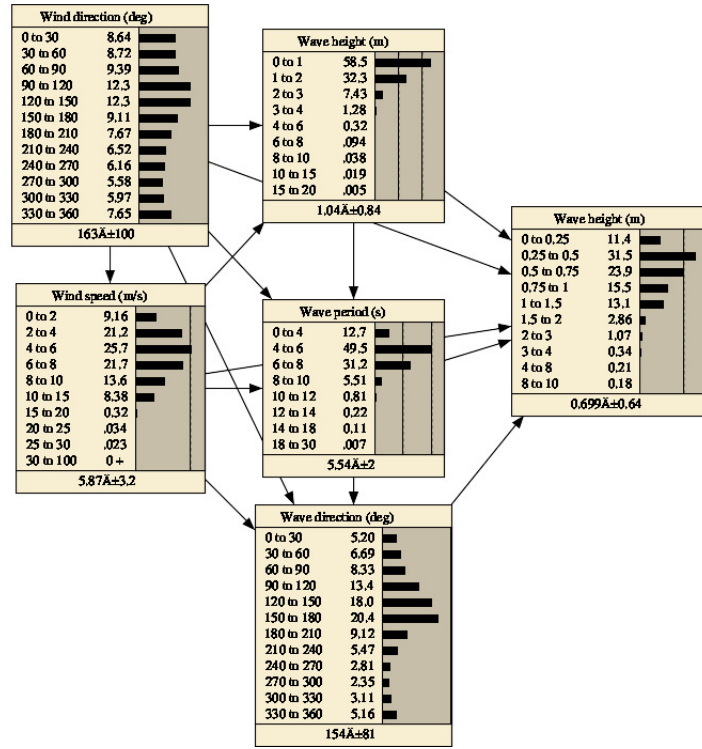


Figure 3: BN for the wave height prediction model.

324 We explore the level of overfitting with a simplified ocean-wave prediction
 325 model based on a BN. The specific BN, illustrated in figure 3, has been used to
 326 drive subsequent predictions of morphologic evolution of a man-made sand berm
 327 constructed near the Chandeleur Islands [21]. Here, we simplify the original
 328 model, which included information from two wave buoys, one tide gage, and
 329 Monte Carlo simulation of a wave-runup model [22]. The first two columns
 330 of network nodes (figure 3) correspond to observations from an offshore buoy
 331 (NOAA 42040) collected from 1996 to 2011. The variables are wind speed
 332 and direction and wave height, period, and direction. The third column has
 333 just one variable—wave height—that was observed at a nearshore buoy between
 334 2000 and 2008. The nearshore buoy was subsequently lost. The BN describing
 335 the prior probabilities of each variable and the conditional probabilities among
 336 variables was designed to resolve boundary conditions at the offshore location
 337 and the prediction at the nearshore location accurately enough to support the
 338 morphologic evolution application. While this BN has very good hindcast skill
 339 (about 0.8), it is not clear that it has equally good forecast skill and whether
 340 fewer probability bins could be retained to give a skillful prediction with optimal
 341 numerical efficiency.

342 Our calibration/validation skill analysis was applied to this net by varying
 343 the number of bins in all variables except for the wave heights. We chose to
 344 resolve these variables consistently with the original model to ensure that prob-
 345 ability predictions spanned a wide range conditions, rather than focusing on the
 346 most probable but extremely low wave-height range that was most common (i.e.,
 347 1-3 m). The number of bins ranged from 2 to 10 for the remaining variables. The
 348 calibration (i.e., hindcast) skill increased for all choices of bin numbers (figure
 349 4). However, the validation skill, averaged over 5 folds, reached its peak value at
 350 4 bins and then decreased dramatically after 6 bins. The optimal bin resolution

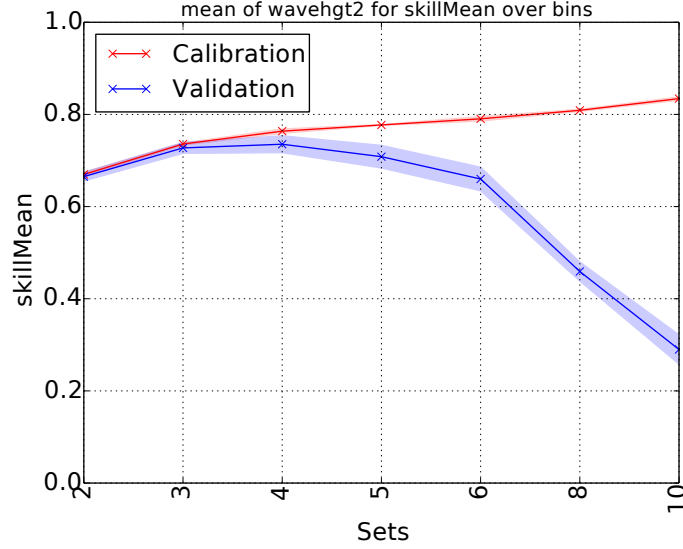


Figure 4: Calibration (descriptive) and Validation (predictive) skill values for various arrangements of bins on the wave prediction BN. For sets identified by a single number, that number of bins was used to discretize all nodes. The color shading indicates the 95% credible interval based on adding and subtracting 2σ to and from the median value of each metric over the 10 folds evaluated.

likely varied for each variable type (wind speed, wind and wave directions, wave
period) and this may explain the flattening of the validation curve between 4
and 6 bins, as it is possible that increasing bin resolution was advantageous,
adding necessary resolution, for some variables but a disadvantage for others.
The rapid decline after 6 bins suggests that none of the variables needed to be
better resolved past this point.

5.2. Model emulation source of groundwater to wells

In the Great Lakes Region of the United States, understanding the interac-
tions between groundwater and surface water are important inputs to ecological
management interests. Specifically, as extraction wells are installed, the base-
flow in streams can be reduced and these reductions can affect fish habitat and
associated societal and economic concerns [23, 24, 25]. An efficient method to

362 determine the source of water to wells has the potential to improve manage-
 363 ment in the region by quickly screening proposed wells. If the source of water
 364 emanating from surface water (either through diversion or depletion) reaches
 365 a management threshold, then further management actions may be triggered.
 366 Using a numerical groundwater model, managers could conceivably explicitly
 367 assess the impact of each proposed well location. But computational run times
 368 and technical background may be prohibitive for that task. A more efficient op-
 369 tion is model emulation as performed on a groundwater model by Fienen et al.
 370 [9].

371 In this case—using MODFLOW-USG [26]—extraction wells were simulated on
 372 multiple staggered grids at sufficient distances that they would not interact in
 373 individual model runs. A base case was also simulated without extraction wells
 374 and, through superposition, the sources of water to the wells was evaluated and
 375 mappable characteristics of each well location were used to create the BN in
 376 Figure 1. For further discussion of the model used in this work see Feinstein
 377 et al. [27].

378 An important question—similar to that evaluated by Fienen et al. [9]—is what
 379 level of complexity provides the best tradeoff between descriptive and predictive
 380 power of the BN. Using CVNetica, it was possible to quickly evaluate k -fold
 381 cross validation for a variety of combinations of bins in the node arrangement
 382 depicted in Figure 1. For the most important response variable—SW_SRC,
 383 which is surface water source—Figure 5 depicts the change in both calibration
 384 and validation performance for 10-fold cross validation performed over an in-
 385 creasingly complex set of bin configurations. While increasing complexity (e.g.,
 386 number of bins per node) monotonically improves calibration (description) over
 387 the training set, the skill improves at first for validation (prediction) but then
 388 degrades dramatically after four bins on the input nodes. In Figure 5, sets iden-

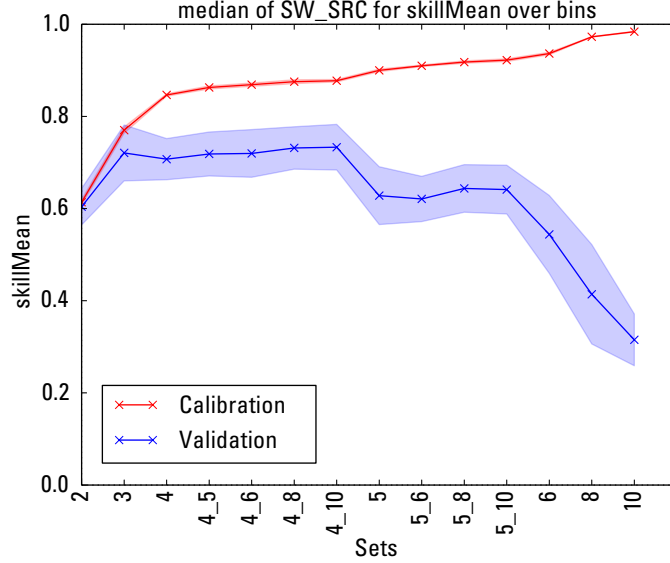


Figure 5: Calibration (descriptive) and Validation (predictive) skill values for various arrangements of bins on the glacial aquifer BN. For sets identified by a single number, that number of bins was used to discretize all nodes. For sets identified by two numbers separated by an underscore, the first and second numbers indicate the number of bins the input and output nodes, respectively, were discretized. The color shading indicates the 95% credible interval based on adding and subtracting 2σ to and from the median value of each metric over the 10 folds evaluated.

389 tified by a single value indicate the number of bins used to discretize each node.
 390 When a second number is present, the first number indicates bins used for the
 391 input nodes while the second indicates bins used for the output nodes. Little
 392 degradation and possibly a slight improvement in validation skill is seen with
 393 increasing output bins for a given complexity of input. This highlights that the
 394 main fitting of the BN takes place with respect to input and output complexity
 395 is more a matter of convenience than a source of real BN complexity.

6. Discussion and Conclusions

396 CVNetica is an open-source Python module using the `ctypes` module to
 397 drive APIs for the Netica BN software. The purpose is to implement cross-
 398 validation techniques for evaluating descriptive (calibration) versus predictive

399 (validation) performance of BNs.

400 We show that CV provides an objective method for determining when a BN
401 is being overfit to the data. And, it appears that overfitting is likely when the
402 BN is designed to resolve physical processes, either observed or modeled. In the
403 cases presented here, the BN design was guided by distribution of the priors
404 of each variable as well as by the intended application of the BN predictions.
405 For instance, in the ocean wave example, the intended application focused on
406 resolving large storm events that were likely to cause erosion. This guided a
407 a choice of fairly high resolution of the input data at the offshore buoy. While
408 physically consistent, a price needed to be paid for the over-resolved output
409 in order to achieve true predictive skill. That price was to greatly reduce the
410 resolution of the input variables from as many as 12 bins to no more than 6
411 bins. While a price is paid in terms of input detail, there is a numerical as well
412 as statistical benefit to reducing the bin resolution. For instance, the original
413 ocean wave BN maintained over a million possible combinations of inputs and
414 outputs scenarios within its conditional probability tables while the optimal 4-
415 bin BN maintained 44 times fewer scenarios, reduced memory requirements, and
416 had increased training and prediction speeds. For instance, the CV processing
417 took 25 minutes for the 4-bin net compared to over 8 hours for the original net
418 (a factor of 20 difference).

419 In the model emulation case, fewer input bins were supported while main-
420 taining good predictive power. Four input bins resulted in good performance
421 while a degradation of predictive skill started with five input bins. Predictive
422 skill was relatively consistent with respect to output bins between 5 and 10 for
423 a given set of input bins. This allows a resource manager to convey outcomes
424 with some flexibility beyond the level of complexity supported by the data on
425 the input side.

426 CVNetica is available for download at (https://github.com/mnfien/NETICA_CV_GENERAL)
427 and the authors welcome proposed contributions to code development going for-
428 ward. These diagnostics and others have the potential to improve the validity
429 of BNs used for prediction in natural resources and other applications.

7. Acknowledgements

This work was funded by (ASIS and GLAS others depending on which example application we use) and the USGS Coastal and Marine Geology Program. We are deeply grateful to Steven Mascaro for his initial PyNetica.py code which he kindly shared as a starting point for this work.

8. Disclaimer

Any use of trade, product, or firm names is for descriptive purposes only and does not imply endorsement by the U.S. Government.

9. References

- [1] F. V. Jensen, T. D. Nielsen, Bayesian networks and decision graphs, Statistics for engineering and information science, Springer, New York, 2001.
- [2] Norsys Software Corp., Netica Version 5.12, <http://www.norsys.com/>, 1990–2013.
- [3] Hugin Expert A/S, HuginExpert Version 2013.0.0, <http://www.hugin.com>, 2013.
- [4] F. Martin de Santa Olalla, A. Dominguez, F. Ortega, A. Artigao, C. Fabeiro, Bayesian networks in planning a large aquifer in Eastern Mancha, Spain, Environmental Modelling and Software 22 (8) (2007) 1089–1100, doi:10.1016/j.envsoft.2006.05.020.

- [5] J. L. Molina, J. Bromley, J. L. García-Aróstegui, C. Sullivan, J. Benavente, Integrated water resources management of overexploited hydrogeological systems using Object-Oriented Bayesian Networks, *Environmental Modelling & Software* 25 (4) (2010) 383–397, doi:10.1016/j.envsoft.2009.10.007.
- [6] J.-L. Molina, D. Pulido-Velázquez, J. L. García-Aróstegui, M. Pulido-Velázquez, Dynamic Bayesian Networks as a Decision Support tool for assessing Climate Change impacts on highly stressed groundwater systems, *Journal of Hydrology* 479 (2013) 113–129, doi:10.1016/j.jhydrol.2012.11.038.
- [7] N. G. Plant, K. T. Holland, Prediction and assimilation of surf-zone processes using a Bayesian network Part I: Forward models, *Coastal Engineering* 58 (1) (2011) 119–130, doi:10.1016/j.coastaleng.2010.09.003.
- [8] N. G. Plant, K. T. Holland, Prediction and assimilation of surf-zone processes using a Bayesian network Part II: Inverse models, *Coastal Engineering* 58 (3) (2011) 256–266, doi:10.1016/j.coastaleng.2010.11.002.
- [9] M. N. Fienen, J. P. Masterson, N. G. Plant, B. T. Gutierrez, E. R. Thieler, Bridging groundwater models and decision support with a Bayesian network, *Water Resources Research* 49 (10) (2013) 6459–6473, ISSN 0043-1397, doi:10.1002/wrcr.20496.
- [10] S. H. Chen, C. A. Pollino, Good practice in Bayesian network modelling, *Environmental Modelling & Software* 37 (2012) 134–145, doi:10.1016/j.envsoft.2012.03.012.
- [11] B. G. Marcot, Metrics for evaluating performance and uncertainty of Bayesian network models, *Ecological Modelling* 230 (2012) 50–62, doi:10.1016/j.ecolmodel.2012.01.013.

- [12] G. Rossum, Python reference manual, Tech. Rep., CWI (Centre for Mathematics and Computer Science), Amsterdam, The Netherlands, The Netherlands, URL <http://www.python.org>, 1995.
- [13] E. Jones, T. Oliphant, P. Peterson, SciPy: Open source scientific tools for Python, URL <http://www.scipy.org/>, 2001–2014.
- [14] K. B. Korb, A. E. Nicholson, Bayesian artificial intelligence, Chapman & Hall/CRC series in computer science and data analysis, Chapman & Hall/CRC, Boca Raton, Fla., 2004.
- [15] E. T. Jaynes, G. L. Bretthorst, Probability theory : the logic of science, Cambridge University Press, Cambridge, UK ; New York, NY, 2003.
- [16] A. Dempster, N. Laird, D. Rubin, Maximum Likelihood From Incomplete Data Via EM Algorithm, Journal of the Royal Statistical Society Series B-Methodological 39 (1) (1977) 1–38.
- [17] Norsys Software Corp., Netica API Programmer’s Library: C Language Version Reference Manual Version 4.18, <http://www.norsys.com/>, 1990–2010.
- [18] T. Hastie, R. Tibshirani, J. H. Friedman, The elements of statistical learning : data mining, inference, and prediction, Springer series in statistics, Springer, New York, 2nd edn., 2009.
- [19] B. T. Gutierrez, N. G. Plant, E. R. Thieler, A Bayesian network to predict coastal vulnerability to sea level rise, Journal of Geophysical Research 116 (F2), doi:10.1029/2010jf001891.
- [20] A. S. Weigend, R. J. Bhansali, PARADIGM CHANGE IN PREDICTION, Philosophical Transactions of the Royal Society a–Mathematical

- Physical and Engineering Sciences 348 (1688) (1994) 405–420, doi: 10.1098/rsta.1994.0100.
- [21] N. G. Plant, J. Flocks, H. F. Stockdon, J. W. Long, K. Guy, D. M. Thompson, J. M. Cormier, J. L. M. C. G. Smith, P. S. Dalyander, Predictions of barrier island berm evolution in a time-varying storm climatology, *Journal of Geophysical Research Earth Surface* 119 (2) (2014) 300–316, doi: 10.1002/2013JF002871.
 - [22] H. F. Stockdon, R. A. Holman, P. A. Howd, J. S. A. H., Empirical parameterization of setup, swash, and runup, *Coastal Engineering* 53 (7) (2006) 573–588, doi:10.1016/j.coastaleng.2005.12.005.
 - [23] F. Ruswick, J. Allan, D. Hamilton, P. Seelbach, The Michigan water withdrawal assessment process—Science and collaboration in sustaining renewable natural resources, *Renewable Resources Journal* 26 (4) (2010) 13–18.
 - [24] P. Barlow, S. Leake, Streamflow depletion by wells—Understanding and managing the effects of groundwater pumping on streamflow, Circular 1376, United States Geological Survey, 2012.
 - [25] K. A. Watson, A. S. Mayer, H. W. Reeves, Groundwater Availability as Constrained by Hydrogeology and Environmental Flows, *Groundwater* 52 (2) (2014) 225–238, doi:10.1111/gwat.12050.
 - [26] S. Panday, C. Langevin, R. Niswonger, M. Ibaraki, J. Hughes, MODFLOW-USG version 1: An unstructured grid version of MODFLOW for simulating groundwater flow and tightly coupled processes using a control volume finite-difference formulation, *Techniques and Methods*, book 6, chap. A45, United States Geological Survey, 2013.

- [27] D. T. Feinstein, M. Fienen, C. Langevin, H. W. Reeves, Application of Unstructured MODFLOW to simulate local groundwater/surface-water interactions at the regional scale as basis for a decision-support tool, Groundwater .