Supporting Information

Improved decision making for water lead testing in U.S. child care facilities using machine-learned Bayesian networks

Riley E. Mulhern,^{1*} AJ Kondash,¹ Ed Norman,² Joseph Johnson,¹ Keith Levine,¹ Andrea McWilliams,¹ Melanie Napier,² Frank Weber,¹ Laurie Stella,¹ Erica Wood,¹ Crystal Lee Pow Jackson,¹ Sarah Colley,¹ Jamie Cajka,¹ Jackie MacDonald Gibson,³ Jennifer Hoponick Redmon^{1*}

1 – RTI International, Research Triangle Park, NC, 27709, USA

2 – Environmental Health Section, Division of Public Health, North Carolina Department of Health and Human Services, Raleigh, NC, 27609, USA

3 – North Carolina State University, Department of Civil, Construction, and Environmental Engineering, Raleigh, NC, 27695, USA

*Corresponding authors: Riley E. Mulhern (<u>rmulhern@rti.org</u>) and Jennifer Hoponick Redmon (<u>jredmon@rti.org</u>)

CONTENTS

|--|

Supporting Figures

Fid	ure S	51.	Flow chart	of mod	el develo	nment and	ner	formance evaluation steps	4
<i>' '</i> S	juic s	·	now churc	oj mou		princine unu	pci.		······································

*Figure S2. Histograms of building-wide maximum and 90th percentile lead concentrations across all centers............*5

 Figure S17. Plots comparing the sensitivity improvement and sampling reduction achieved by each model compared to the alternative heuristics. Models predicting the maximum lead level (max >1, >5, >10, and >15) are shown on the left. Models predicting the 90th percentile lead level (P90 >1, >5, >10, >15) are shown on the right...19

Supporting Tables

Table S1. Variables included in the Clean Water for Carolina Kids data set ¹⁶ used as predictors in the BN models predict building-wide lead risk	to 20
Table S2. Additional variables included in the data set compiled from publicly available data sources.	22
Table S3. Summary of all performance metrics for each model	24
Table S4. Summary of mean F1 and F2 scores for each heuristic from 10-fold cross validation	24
Table S5. Summary table of all significant variables included in each model. Interactive model structures can also seen at www.cleanwaterforcarolinakids.org/publications/bn models.	o be 25



Figure S1. Flow chart of model development and performance evaluation steps.



Figure S2. Histograms of building-wide maximum and 90th percentile lead concentrations across all centers.

Additional Methods Information

Discretization

Several other supervised discretization algorithms were also tested, such as Minimum Description Length Principle (MDLP), Class-Attribute Interdependence Maximization (CAIM), and Ameva algorithms (each of which has been shown to adequately represent the true data distribution with a minimal number of discrete intervals^{1–3}), but these were found to inadvertently remove certain variables before the machine learning step by failing to select any cut points, thus rendering them unusable in a BN model. Since variables in both the training and test set must have the same intervals for BN models, the cut points selected by the discretization algorithm for the training set were then manually applied to the test set. Although discretization of the complete data set prior to machine learning (rather than performing separate discretization steps on the training and test sets) is common in BN modelling,^{4–6} the manual application of the training set cut points to the test set minimizes potential data leakage concerns.

Structure learning and feature selection

Additional measures of arc strength were assessed alongside Chi-squared tests, including whether removing the node from the network would result in a reduction of the model's Akaike Information Criterion or Bayesian Information Criterion, and whether the node was connected to the target at least 10% of the time during 100 bootstrapped unsupervised structure learning iterations. The various variable selection approaches performed comparably with the exception of the Bayesian Information Criterion (**Figure S3**). Chi-squared tests were chosen because they allowed missing values in the training set to be handled transparently, whereas the other approaches required imputation as a pre-processing step (see "Missing values processing" below). An example of the arc strengths calculated by the Chi-squared tests can be seen for the Max≥1 model in **Figure S4**. The bootstrapping procedure yielded significantly fewer variables and could be used to simplify the models further if many of the variables selected by Chi-squared tests were unavailable in certain areas.



Figure S3. Comparison of model performance after 10-fold cross validation using various variable selection criteria, including AIC, BIC, Chi-squared tests, and bootstrapping. The various criteria yielded similarly performing models except for BIC. For this comparison, missing values in the training set were imputed using structural expectation maximization. Chi-squared tests were ultimately selected as the variable selection criteria for their ability to handle missing values.



Figure S4. Strength of association between the target and all possible predictor nodes for model Max>=1 determined by Chi-Squared tests. Predictors not directly associated with the target node (dashed lines) were removed from the model.

Missing values processing

The final data set contained only 6.2% missing values. While there is no recognized cutoff for an acceptable level of missingness for statistical analysis, <10% missing data is generally considered low, with little bias expected on the result.⁷ However, a little over half of the facilities were missing data for at least one field. The missingness pattern can be seen in **Figure S5**.



Figure S5. Missingness pattern of the complete data set. Overall, the data set contained only 6.2% missing values, but over half of the facilities were missing data for at least one field.

Most of the missing data were related to enrollment and demographic information (e.g., percent non-white students, percent free/reduced lunch students, total enrollment) that was not entered or was entered incorrectly by the facility during their registration with the Clean Water for Carolina Kids program (for example, erroneous enrollment numbers resulting in >100% non-White; in these cases, the values were treated as missing). These data are considered missing not at random (MNAR) since the missingness mechanism is likely a function of the variable itself (e.g., centers with a higher proportion of children of color may be more likely to accidently misreport their student demographic information compared to centers with only White students, or centers with a higher proportion of free/reduced lunch children may not want this information to be known due to stigma around receiving free services⁸). Additional missing data were related to Lead and Copper Rule monitoring and water system information in cases where facility staff either did not know or did not report which community water system they were served by, were served by a private well, or operated as a NTNC water system. We consider these data to be missing at random (MAR) because their missingness depends on other factors in the data set, such as water system type.

To assess the effect of missingness in the training data set on model development, we conducted a sensitivity analysis to evaluate several different approaches to handling missing values,⁹ including multiple imputation by chained equations (MICE),^{10,11} structural expectation maximization (SEM) from a separate unsupervised BN,^{11,12} estimating conditional probability distributions for each network node using complete observations for that node's local network.¹³ The results of this sensitivity analysis can be seen in **Figure S6**. Overall, the MI and SEM imputation approaches provided negligible benefit in performance while significantly increasing computational costs compared to handling missing values transparently.

Meanwhile, missing values could not be ignored in the test data set since the "predict" function in *bnlearn* requires complete evidence to predict the probability of the target node. One advantage of BNs is that, once trained, the network structure itself can be used to estimate missing values in a test set using Bayesian likelihood weighting.¹⁴ Given the potential difficulty of collecting building-specific data in some cases, missing values are highly likely in practice. Thus, practitioners could feasibly apply our models to new data without having to develop a separate imputation procedure to handle missing values before prediction. To eliminate data leakage concerns and to best simulate the model's practical performance where no prior knowledge of the outcome is available, our model strips the test data (whether the external test set, the test fold during k-fold cross validation, or the training data itself when assessing the internal model performance) of the real result of the target outcome and generates

S9

a random guess in its place. A random guess is necessary because the "impute" function in *bnlearn* needs a data set that matches the structure of the network, so the target variable cannot be blank. Importantly, this simulates how practitioners would approach the model with missing data and does not introduce any "unfair" information to the model during prediction. The random guess is then removed after the imputation step and the predicted probability of the target node is calculated using the imputed testing data as evidence.



Figure S6. Comparison of model performance after 10-fold cross validation using various approaches to handle missing values in the data set, including using complete observations from each nodes local network, multiple imputation by chained equations (MICE) and structural expectation maximization (SEM). Overall, the results were not significantly affected by the imputation procedure.



Figure S7. Structures of the eight final models. Interactive versions can also be seen at: www.cleanwaterforcarolinakids.org/publications/bn_models



Figure S8. Receiver operating characteristic (ROC) curves for each model. Models predicting the maximum lead level (max >1, >5, >10, and >15) are shown on the left. Models predicting the 90th percentile lead level (P90 >1, >5, >10, >15) are shown on the right.



Figure S9. Relationship between the prior probability of each model, indicating the level of class imbalance, and the AU-ROC (pink), AU-PR (grey), and F2-score (gold) performance metrics (F1-scores are not shown for visual clarity, but follow a similar trend as F2-scores). AU-ROC values are insensitive to large class imbalance, while AU-PR and F2-scores reveal decreasing predictive performance for more rare outcomes.



Figure S10. Precision-recall (PR) curves for each model. Models predicting the maximum lead level (max >1, >5, >10, and >15) are shown on the left. Models predicting the 90^{th} percentile lead level (P90 >1, >5, >10, >15) are shown on the right.



Figure S11. Posterior probabilities of the four most frequently selected variables across all eight models compared to each model's prior. The width of the bars shows the range of impact of each variable on building-wide lead risk. Grey diamonds indicate the prior probability of each model target considering all states. Models are ordered on the y-axis according to the size of the variable's impact on the target, with models where the variable was more impactful at top.



Figure S12. Posterior probabilities for the variable "Past Faucet Fixture Change" for the seven models where it was included. If the center did not know whether any past faucets had been changed significantly increased the building-wide water lead risk.



Figure S13. Posterior probabilities for the variables "Home-based" (Panel A) and "School-based" (Panel B) highlighting decreases in risk for small home-based centers, and increases in risk for larger school-based centers for certain model targets.



Figure S14. Posterior probabilities for the variables "Phosphate addition" (Panel A) and "pH adjustment" (Panel B) highlighting increases in water lead risk in child care centers served by water systems (public or private) that do not use phosphate-based corrosion inhibitors or pH adjustment.



Figure S15. Posterior probabilities for the variable "Past LCR exceedance" for the five models where it was included. If the center relied on a private well it was coded as "LCR NA," indicating that Lead and Copper Rule monitoring is not applicable. Private well water significantly increased the building-wide water lead risk, whereas a past action level (AL) exceedance by the water utility decreased the risk.



Figure S16. Tornado charts showing the minimum and maximum posterior probabilities and node states for all variables included in each model. The dashed vertical line in each chart shows the prior probability of the target for each model. Posterior probabilities to the right indicate increased risk; posterior probabilities to the left indicate decreased risk. The width of the line indicates the magnitude of the effect of the variable on the target.



Figure S17. Plots comparing the sensitivity improvement and sampling reduction achieved by each model compared to the alternative heuristics. Models predicting the maximum lead level (max >1, >5, >10, and >15) are shown on the left. Models predicting the 90th percentile lead level (P90 >1, >5, >10, >15) are shown on the right.

Table S1. Variables included in the Clean Water for Carolina Kids data set¹⁵ used as predictors in the BN models to predict building-wide lead risk.

Variable	Abbreviated name	Description	Variable type	Summary	
Name		-		-	
Building information	tion	ł	•		
Ownership	OWN_OR_LEASE	Whether the center	Binary	Own: 2,774	
		building was owned or		Lease: 1,229	
		leased.			
School-based	school	Whether the center	Binary	School-based: 742	
		was located in a		Not school-based:	
I I a sea a la sea d	have have	school building	D'a sur	3,261	
Home-based	nome_based	whether the center	Binary	Home-based: 192	
		was located in a nome		3 811	
Franchised	franchised	Whether the center	Binary	Franchised: 205	
Tuneniscu	maneniseu	was part of a franchise	Dinary	Not franchised: 3.798	
Year built	BUILT cat	Year the building was	Categorical	Pre-1988: 1.944	
		constructed	(Pre-1988, 1988-2013,	1988-2013: 1,875	
			or Post 2014)	After 2014: 184	
Community water	CWS	Whether the center	Binary	Connected to CWS:	
system		was connected to a		3,452	
		public community		Not connected to	
		water system (CWS)		CWS: 551	
# samples	nsamples	The number of	Discrete	Mean: 5.7	
		drinking and cooking		Max: E1	
		collected for analysis		Missing: 0	
		by the center (center		wissing. U	
		staff were instructed			
		to sample all drinking			
		and cooking taps).			
Percent filtered	perc filtered	The percentage of	Continuous	Mean: 12%	
		taps sampled by each		Median: 0%	
		center that had a		Max: 100%	
		point-of-use filter		Missing: 96 (2.4%)	
		installed and was			
		flagged as filtered by			
		center staff.			
Private well	private_well	Whether the center	Binary	Private well: 175	
		used a private well for		No private well: 3,828	
		its water supply. This			
		was determined			
		the conter responded			
		"No" to the question			
		"Does your drinking			
		water come from a			
		public water			
		treatment plant?" and			
		"Yes" to the question			
		"Does the center use			
		well water?" This			
		variable was not			
		mutually exclusive			
		with "Community			
1		water system"	1		

6)
•/
82
e:

Variable	Abbreviated	Description	Variable	Summary	Data	
Name	name	-	type	-	source	
Past LCR exceedance	any_lcr_exceedance	Whether the facility was served by a water system with an action level exceedance (≥10% of collected water samples collected by the utility exceeding 15 ppb) in the five years prior to Clean Water for Carolina Kids sampling (2016-2021). If served by a private well, coded as "LCR NA" but not missing.	Categorical	Past exceedance: 55 No past exceedance: 2,889 LCR NA: 4.3% Missing: 884 (22%)	EPA Safe Drinking Water Information System ¹⁶	
type	type_binary	water system source water type. Groundwater use indicates groundwater alone; if a water system used any surface water it was coded as surface water. Centers reporting using private well water default to groundwater. Other centers not served by a community water system coded as missing.	Binary	Surface water: 2,772 Missing: 622 (15.5%)	North Carolina Drinking Water Watch ¹⁷	
# connections of water system	connections_cat	Number of service connections in the water system. Centers reporting using private well water default to 1 service connection. Other centers not served by a community water system coded as missing.	Categorical	<3300 (Small, very small, and private systems): 766 3301-10000 (Medium systems): 744 >10001 (Large, very large systems): 1971 Missing: 622 (15.5%)	North Carolina Drinking Water Watch ¹⁷	
Phosphate addition	Phos_binary	Whether the water system implements phosphate-based corrosion control. Centers reporting using private well water default to no phosphate addition. Other centers not served by a community water	Binary	Phosphate addition: 1,944 No phosphate addition: 1,438 Missing: 622 (15.5%)	North Carolina Drinking Water Watch ¹⁷	

 Table S2. Additional variables included in the data set compiled from publicly available data sources.

		system coded as			
		missing.			
pH adjustment	pH_binary	Whether the water	Binary	pH adjustment:	North
		system implements		2,481	Carolina
		pH adjustment.		No pH adjustment:	Drinking
		Centers reporting		901	Water
		using private well		Missing: 622 (15.5%)	Watch ¹⁷
		water default to no			
		pH adjustment. Other			
		centers not served by			
		a community water			
		system coded as			
		missing.			
Coagulation	coagulation	Whether the water	Binary	Coagulation: 1,947	North
		system implements		No coagulation:	Carolina
		coagulation in the		1,435	Drinking
		treatment train.		Missing: 622 (15.5%)	Water
		Centers reporting			Watch ¹⁷
		using private well			
		water default to no			
		coagulation. Other			
		centers not served by			
		a community water			
		system coded as			
		missing.			
Chloramines	chloramines	Whether the water	Binary	Chloramines: 1,032	North
		system uses		No chloramines:	Carolina
		chloramination for		2,350	Drinking
		disinfection. Centers		Missing: 622 (15.5%)	Water
		reporting using			Watch ¹⁷
		private well water			
		default to no			
		chloramination.			
		Other centers not			
		served by a			
		community water			
		system coded as			
		missing.			
Urbanicity	ruca_cat	Whether the center	Categorical	Metropolitan: 2,920	US
		was located in a		Micropolitan: 669	Department
		metropolitan,		Rural: 161	of
		micropolitan, small		Small town: 24	Agriculture,
		town, or rural area		Missing: 8 (0.1%)	Economic
		based on Rural-Urban			Research
		Commuting Area			Service ¹⁸
		(RUCA) code			
Block group	med_hh_income_cbg	Median household	Continuous	Mean: 56,031	American
median		income of the block		Median: 49,458	Community
household		group where each		Max: 250,001	Survey
income		facility was located		Missing: 174 (4.3%)	2020 ¹⁹
Block group	perc_hs_higher_cbg	Proportion of block	Continuous	Mean: 0.13	American
educational		group population		Median: 0.11	Community
attainment		having attained a high		Max: 0.77	Survey
		school degree or		Missing: 11 (0.2%)	2020 ¹⁹
		higher			
Block group %	perc_non_white_cbg	Proportion of block	Continuous	Mean: 0.39	American
non-White		group population that		Median: 0.34	Community
		identified as a race		Max: 1.0	

category other than	Missing: 11 (0.2%)	Survey
White alone		202019

 Table S3. Summary of all performance metrics for each model.

	Prior probability		Mean AU-ROC 10-fold cross validation on	AU- ROC		Mean AU-PR 10-fold cross validation on	AU- PR	Мах	Мах
Model	of model	AU-ROC full	training	test	AU-PR	training	test	F1-	F2-
name	target	training set	set	set	training set	set	set	score	score
Max>1	0.56	0.75	0.73	0.70	0.78	0.76	0.73	0.76	0.87
Max>5	0.21	0.72	0.71	0.76	0.43	0.43	0.38	0.48	0.62
Max>10	0.13	0.75	0.72	0.74	0.32	0.31	0.26	0.38	0.52
Max>15	0.09	0.73	0.69	0.76	0.24	0.24	0.20	0.31	0.43
P90>1	0.49	0.71	0.68	0.66	0.70	0.68	0.63	0.69	0.84
P90>5	0.15	0.71	0.67	0.71	0.31	0.31	0.22	0.37	0.52
P90>10	0.07	0.72	0.65	0.75	0.21	0.20	0.14	0.25	0.37
P90>15	0.05	0.67	0.62	0.69	0.12	0.16	0.10	0.22	0.28

 Table S4. Summary of mean F1 and F2 scores for each heuristic from 10-fold cross validation.

		F1 sc	ores	F2 scores			
Heuristic	Mean	Range	BN model mean %	Mean	Range	BN model mean %	
			improvement			improvement	
Groundwater	0.23	0.15-0.31	80%	0.23	0.20-0.25	142%	
Head Start	0.23	0.14-0.28	83%	0.22	0.19-0.25	155%	
Pre-1988	0.00	0.11-0.53	65%	0.36	0.22-0.53	51%	
buildings	0.28						
Private well	0.14	0.10-0.19	248%	0.11	0.07-0.16	506%	
water	0.14						

Model	Variables included	Model	Variables included
	% free/reduced lunch enrollment		% free/reduced lunch enrollment
	% non-White enrollment		# samples
	Total enrollment		% taps filtered
	# samples		Head Start
	% taps filtered	1	School-based
	Head Start		Home-based
	School-based		Past faucet fixture change
Max>1	Home-based	P90>1	Year of past faucet fixture change
IVIAX-1	Past faucet fixture change		Year center began operating
	Year of past faucet fixture change		Source water type
	Year center began operating		pH adjustment
	Source water type		Chloramination
	pH adjustment		# connections of water system
	Chloramination		Urbanicity
	# connections of water system		Block group median household income
	Urbanicity		% free/reduced lunch enrollment
	% free/reduced lunch enrollment		# samples
	Total enrollment		Head Start
	# samples		School-based
	Head Start		Home-based
Max>5	School-based	P90>5	Past faucet fixture change
WIGAP 5	Home-based	-	Source water type
	Past faucet fixture change		# connections of water system
	Year center began operating		Past LCR exceedance
	Source water type		Block group % non-White
	# connections of water system		On-site wastewater system
	% free/reduced lunch enrollment		% free/reduced lunch enrollment
	# samples		# samples
	Head Start		Head Start
	School-based		School-based
Max>10	Home-based		Home-based
	Past faucet fixture change	P90>10	Past faucet fixture change
	Source water type	1 30/ 10	Source water type
	# connections of water system		pH adjustment
	Past LCR exceedance		Past LCR exceedance
	Phosphate addition		Phosphate addition
Max>1F	% free/reduced lunch enrollment		Block group % non-White
IVIAX/13	# samples		On-site wastewater system

 Table S5.
 Summary table of all significant variables included in each model. Interactive model structures can also be seen at www.cleanwaterforcarolinakids.org/publications/bn models.

Head Start		Block group educational attainment		
School-based		% free/reduced lunch enrollment		
Home-based		# samples		
Past faucet fixture change	P90>15	Head Start		
Year center began operating		Source water type		
Source water type		pH adjustment		
pH adjustment		Past LCR exceedance		
Past LCR exceedance		On-site wastewater system		
Phosphate addition				
Block group % non-White				

References

- Gonzalez-Abril, L.; Cuberos, F. J.; Velasco, F.; Ortega, J. A. Ameva: An Autonomous Discretization Algorithm. *Expert Syst Appl* 2009, *36* (3 PART 1), 5327–5332. https://doi.org/10.1016/j.eswa.2008.06.063.
- Ropero, R. F.; Renooij, S.; van der Gaag, L. C. Discretizing Environmental Data for Learning Bayesian-Network Classifiers. *Ecol Modell* 2018, *368*, 391–403. https://doi.org/10.1016/j.ecolmodel.2017.12.015.
- García, S.; Luengo, J.; Sáez, J. A.; López, V.; Herrera, F. A Survey of Discretization Techniques: Taxonomy and Empirical Analysis in Supervised Learning. *IEEE Trans Knowl Data Eng* 2013, 25 (4), 734–750. https://doi.org/10.1109/TKDE.2012.35.
- (4) Conrady, S.; Jouffe, L. *Bayesian Networks and BayesiaLab: A Practical Introduction for Researchers*; Bayesia USA: Franklin, TN, 2015.
- Roostaei, J.; Colley, S.; Mulhern, R.; May, A. A.; Gibson, J. M. Predicting the Risk of GenX Contamination in Private Well Water Using a Machine-Learned Bayesian Network Model. J Hazard Mater 2021, 411 (October 2020), 125075. https://doi.org/10.1016/j.jhazmat.2021.125075.
- (6) Mulhern, R.; Roostaei, J.; Schwetschenau, S.; Pruthi, T.; Campbell, C.; MacDonald Gibson, J. A New Approach to a Legacy Concern: Evaluating Machine-Learned Bayesian Networks to Predict Childhood Lead Exposure Risk from Community Water Systems. *Environ Res* 2022, 204, 112146. https://doi.org/10.1016/j.envres.2021.112146.
- (7) Dong, Y.; Peng, C. Y. J. Principled Missing Data Methods for Researchers. Springerplus 2013, 2 (222), 1–17. https://doi.org/10.1186/2193-1801-2-222.
- Kuhns, C.; Martinchek, K.; Gupta, P. Combating Food Insecurity and Supporting Child Nutrition through the Child Care Sector; Washington D.C., 2021. https://www.urban.org/sites/default/files/publication/105106/combating-food-insecurity-andsupporting-child-nutrition-through-the-child-care-sector.pdf.
- Madley-Dowd, P.; Hughes, R.; Tilling, K.; Heron, J. The Proportion of Missing Data Should Not Be Used to Guide Decisions on Multiple Imputation. *J Clin Epidemiol* 2019, *110*, 63–73. https://doi.org/10.1016/j.jclinepi.2019.02.016.
- van Buuren, S.; Groothuis-Oudshoorn, K. Mice: Multivariate Imputation by Chained Equations in R. J Stat Softw 2011, 45 (3), 1–67. https://doi.org/10.18637/jss.v045.i03.
- Ke, X.; Keenan, K.; Smith, V. A. Treatment of Missing Data in Bayesian Network Structure Learning: An Application to Linked Biomedical and Social Survey Data. *BMC Med Res Methodol* 2022, 9, 1–2. https://doi.org/10.1186/s12874-022-01781-9.
- (12) Scutari, M. *Structure learning from missing data*. bnlearn an R package for Bayesian network learning and inference. https://www.bnlearn.com/documentation/man/structural.em.html (accessed 2023-02-09).

- (13) Scutari, M. Parameter learning from data with missing values. bnlearn an R package for Bayesian network learning and inference. https://www.bnlearn.com/examples/missing-fitting/ (accessed 2023-02-09).
- Scutari, M. Predict or impute missing data from a Bayesian network. bnlearn an R package for Bayesian network learning and inference. https://www.bnlearn.com/documentation/man/predict.and.impute.html (accessed 2023-02-09).
- Hoponick Redmon, J.; Kondash, A. J.; Norman, E.; Johnson, J.; Levine, K.; McWilliams, A.; Napier, M.; Weber, F.; Stella, L.; Wood, E.; Lee Pow Jackson, C.; Mulhern, R. Lead Levels in Tap Water at Licensed North Carolina Child Care Facilities, 2020-2021. *Am J Public Health* 2022, *112* (S7), S695–S705. https://doi.org/10.2105/AJPH.2022.307003.
- (16) EPA. Safe Drinking Water Information System. https://sdwis.epa.gov/ords/sfdw_pub/f?p=108:200 (accessed 2023-02-07).
- (17) NCDEQ. *Drinking Water Watch*. https://www.pwss.enr.state.nc.us/NCDWW2/ (accessed 2020-07-05).
- USDA Economic Research Service. *Rural-Urban Commuting Area Codes*. https://www.ers.usda.gov/data-products/rural-urban-commuting-area-codes.aspx (accessed 2022-08-15).
- (19) Manson, S.; Schroder, J.; van Riper, D.; Kugler, T.; Ruggles, S. *IPUMS National Historical Geographic Information System: Version 17.0*; Minneapolis, MN, 2022. https://doi.org/10.18128/D050.V17.0.