

Appendix 7. Model Archival Summary for Hardness Concentration at U.S. Geological Survey Site 06888990, Kansas River above Topeka Weir at Topeka, Kansas, during November 2018 through June 2021

This model archival summary summarizes the hardness as calcium carbonate (CaCO_3 ; U.S. Geological Survey [USGS] parameter code 00900) concentration model developed to compute 15-minute CaCO_3 concentrations from November 2018 onward. This model is specific to USGS site 06888990, the Kansas River above Topeka Weir at Topeka, Kansas, during this study period and cannot be applied to data collected from other sites on the Kansas River or data collected from other waterbodies.

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

Site and Model Information

Site number: 06888990

Site name: Kansas River above Topeka Weir at Topeka, Kans.

Location: Lat 39°04'19", long 95°42'58" referenced to North American Datum of 1927, in NW 1/4 sec.23, T.11 S., R.15 E., Shawnee County, Kans., hydrologic unit 10270102.

Equipment: A Xylem YSI EXO2 water-quality monitor equipped with sensors for water temperature, specific conductance (SC), dissolved oxygen, pH, turbidity, and chlorophyll and phycocyanin fluorescence was installed during November 2018 through June 2021. Readings from the water-quality monitor were recorded every 15 minutes and transmitted by way of satellite, hourly.

Date model was created: December 8, 2021

Model-calibration data period: November 28, 2018, through June 21, 2021

Model-application date: November 28, 2018, onward

Model-Calibration Dataset

All data were collected using USGS protocols (Wagner and others, 2006; U.S. Geological Survey, variously dated) and are stored in the USGS National Water Information System (U.S. Geological Survey, 2022) database and available to the public. Ordinary least squares analysis was used to develop regression models using R programming language (R Core Team, 2022). Potential explanatory variables that were evaluated individually and in combination included streamflow, water temperature, SC, dissolved oxygen, pH, turbidity, and chlorophyll and phycocyanin fluorescence. These potential explanatory variables were interpolated within the 15-minute continuous record based on sample time. The maximum time span between two continuous data points used for interpolation was 2 hours (in order to preserve the sample dataset, field monitor averages obtained during sample collection were used for model development data if no continuous data were available or if gaps larger than 2 hours in the continuous data record resulted in missing interpolated data). Seasonal components (sine and cosine variables) also were evaluated as potential explanatory variables. Previously published explanatory variables (Rasmussen and others, 2005; Foster and Graham, 2016; Williams, 2021) at other Kansas River sites were strongly considered for continuity in model form.

The final selected regression model was based on 34 concurrent measurements of CaCO_3 concentration and sensor-measured SC during November 28, 2018, through June 21, 2021. Samples were collected throughout the range of continuously observed hydrologic conditions. No samples had concentrations below laboratory minimum reporting limits.

Potential outliers initially were identified using scatterplots of the CaCO_3 and SC model-calibration data (Rasmussen and others, 2009). Studentized residuals from the model were inspected for values greater than three or less than negative three (Pardoe, 2020). Values outside of that range were considered potential outliers and were investigated. Additionally, computations of leverage, Cook's distance (Cook's D), and difference in fits (DFFITS) statistics were used to estimate potential outlier effect on the final selected regression model (Cook, 1977; Helsel and others, 2020). Outliers were investigated for potential removal from the model-calibration dataset by confirming correct database entry, evaluating laboratory analytical performance, and reviewing field notes associated with the sample in question (Rasmussen and others, 2009). All potential outliers were not determined to have errors associated with sample collection, processing, or analysis and were therefore considered valid.

Hardness Sampling Details

During November 2018 through February 2019, samples were collected using the equal-width increment collection method (U.S. Geological Survey, variously dated). In March 2019, sample collection location changed to the southern bank of the Kansas River above Topeka Weir using the single-vertical collection method (U.S. Geological Survey, variously dated) to avoid safety risks caused by a nearby low-head dam. All samples were composited for analysis (U.S. Geological Survey, variously dated). During

November 2018 through June 2020, samples were collected on a biweekly to monthly basis. During July 2020 through June 2021, samples were collected on a monthly to quarterly basis, depending on flow conditions. Samples occasionally were collected during targeted reservoir release and runoff events to get a more representative dataset. A FISP US DH-81, DH-95, D-95, or D-96a depth integrating sampler was used. Samples were analyzed for CaCO₃ concentration at the USGS National Water Quality Laboratory in Lakewood, Colorado.

Model Development

Ordinary least squares regression analysis was done using the *stats* (v4.3.0) package in R programming language (R Core Team, 2022) to relate discretely collected CaCO₃ concentration to sensor-measured SC. The distribution of residuals (the difference between the measured and computed values) was examined for normality, and the plots of residuals were examined for homoscedasticity (departures from zero did not change substantially over the range of computed values).

SC was selected as a good surrogate for CaCO₃ based on residual plots, coefficient of determination (R^2), and model standard percentage error. Values for all the aforementioned statistics, all relevant sample data, and additional statistical information are included in the Model Statistics, Data, and Plots section of this appendix.

Model Summary

The following is a summary of the final regression analysis for CaCO₃ concentration at USGS site 06888990:

CaCO₃ concentration-based model:

$$\log CaCO_3 = 0.771(\log SC) + 0.201$$

where

log = logarithm base 10,

CaCO₃ = hardness as calcium carbonate concentration, in milligrams per liter, and

SC = specific conductance, in microsiemens per centimeter at 25 degrees Celsius.

SC makes physical and statistical sense as an explanatory variable for CaCO₃ because of its positive correlation with charged ionic species (Hem, 1985).

The logarithmically (log) transformed model may be retransformed to the original units so that CaCO₃ can be calculated directly. The retransformation introduces a bias in the calculated constituent. This bias may be corrected using Duan's bias correction factor (BCF; Duan, 1983). For this model, the calculated BCF is 1.00. The retransformed model, accounting for BCF is as follows:

$$CaCO_3 = 1 \times (SC^{0.771} \times 10^{0.201})$$

This model was developed using continuous and discrete water-quality data collected during November 2018 through June 2021. These data were collected throughout the observed range of streamflow conditions during this time. However, a limitation in model accuracy during conditions outside of those observed during November 2018 through June 2021 should be considered when interpreting model computations beyond June 2021.

Previous Models

There are no previously published models at this site. However, similar models have been published at other Kansas River sites, as documented by Rasmussen and others (2005), Foster and Graham (2016), and Williams (2021).

Model Statistics, Data, and Plots

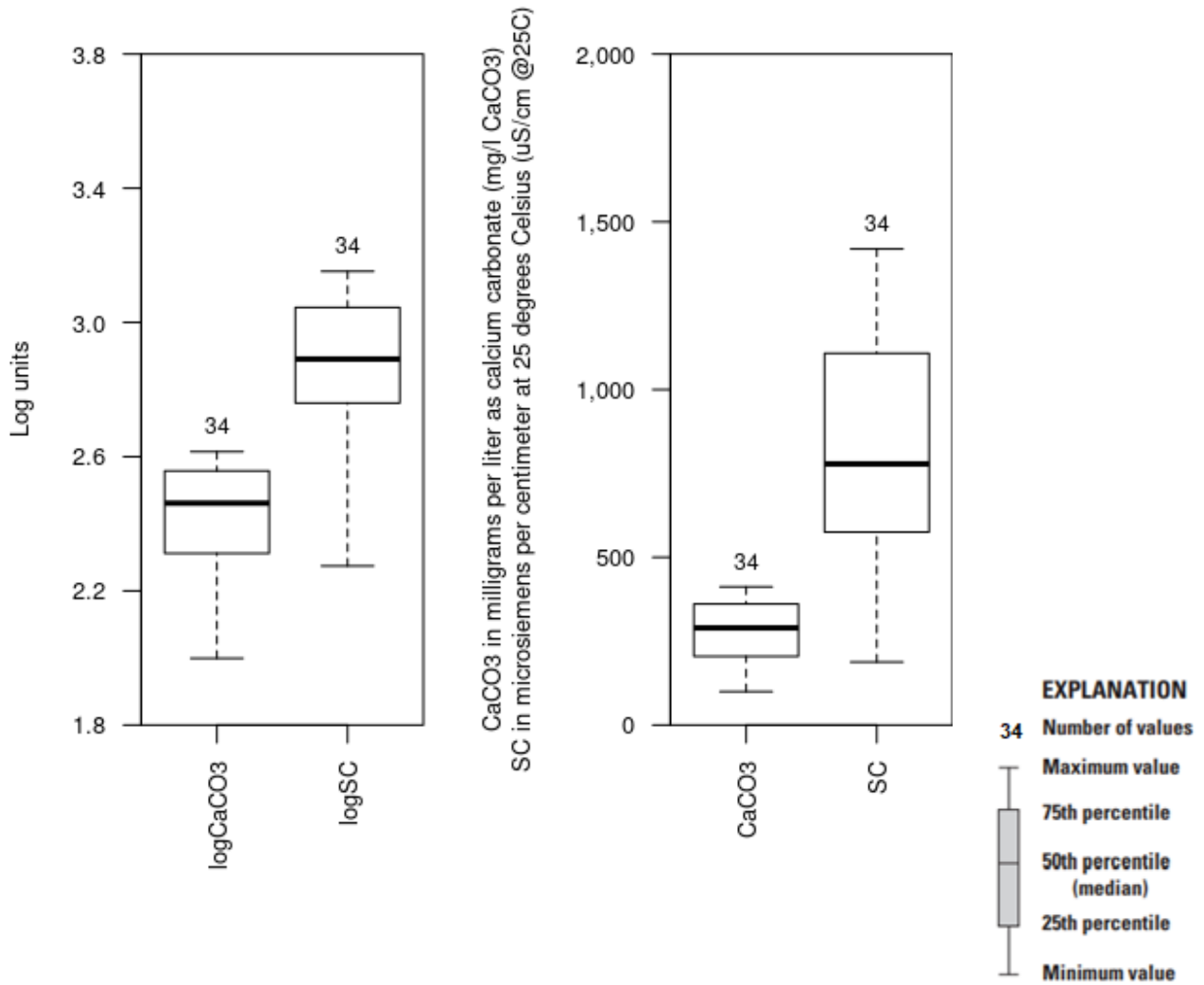
Model

$$\log CaCO_3 = 0.771(\log SC) + 0.201$$

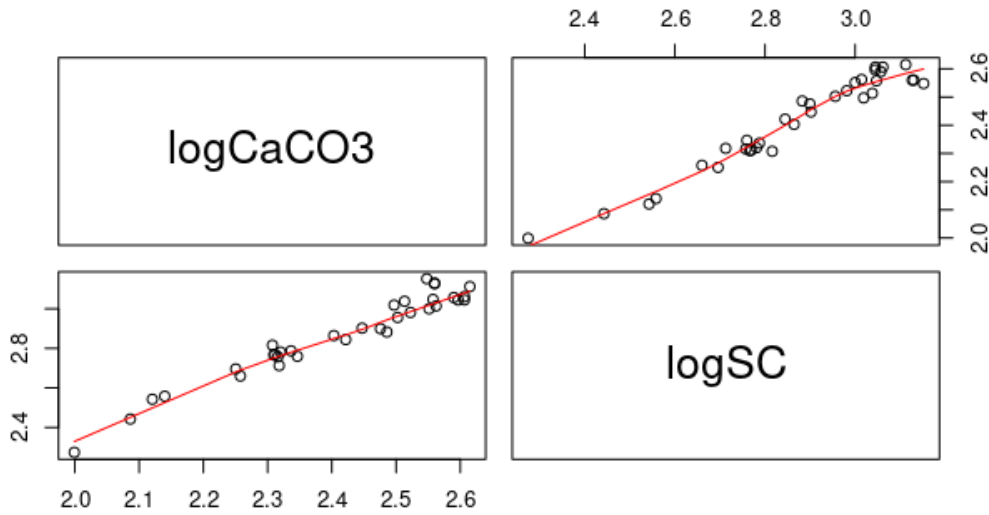
Variable Summary Statistics

	logCaCO3	CaCO3	logSC	SC
Minimum	2.00	99.8	2.27	188
1st Quartile	2.31	205	2.76	575
Median	2.46	289	2.89	778
Mean	2.41	276	2.87	814
3rd Quartile	2.56	361	3.04	1,110
Maximum	2.61	412	3.15	1,420

Box Plots



Exploratory Plots



Red line shows the locally weighted scatterplot smoothing (LOWESS).

The x- and y-axis labels for a given bivariate plot are defined by the intersecting row and column labels.

Basic Model Statistics

Number of observations	34
Standard error (RMSE)	0.0381
Mean model standard percentage error (MSPE)	8.78
Coefficient of determination (R^2)	0.949
Adjusted coefficient of determination (Adj. R^2)	0.947
Bias correction factor (BCF)	1.00

Explanatory Variables

	Coefficients	Standard Error	t value	Pr(> t)
(Intercept)	0.201	0.0909	2.21	3.44e-02
logSC	0.771	0.0316	24.40	3.06e-22

Correlation Matrix

	Intercept	E.vars
Intercept	1.000	-0.997
E.vars	-0.997	1.000

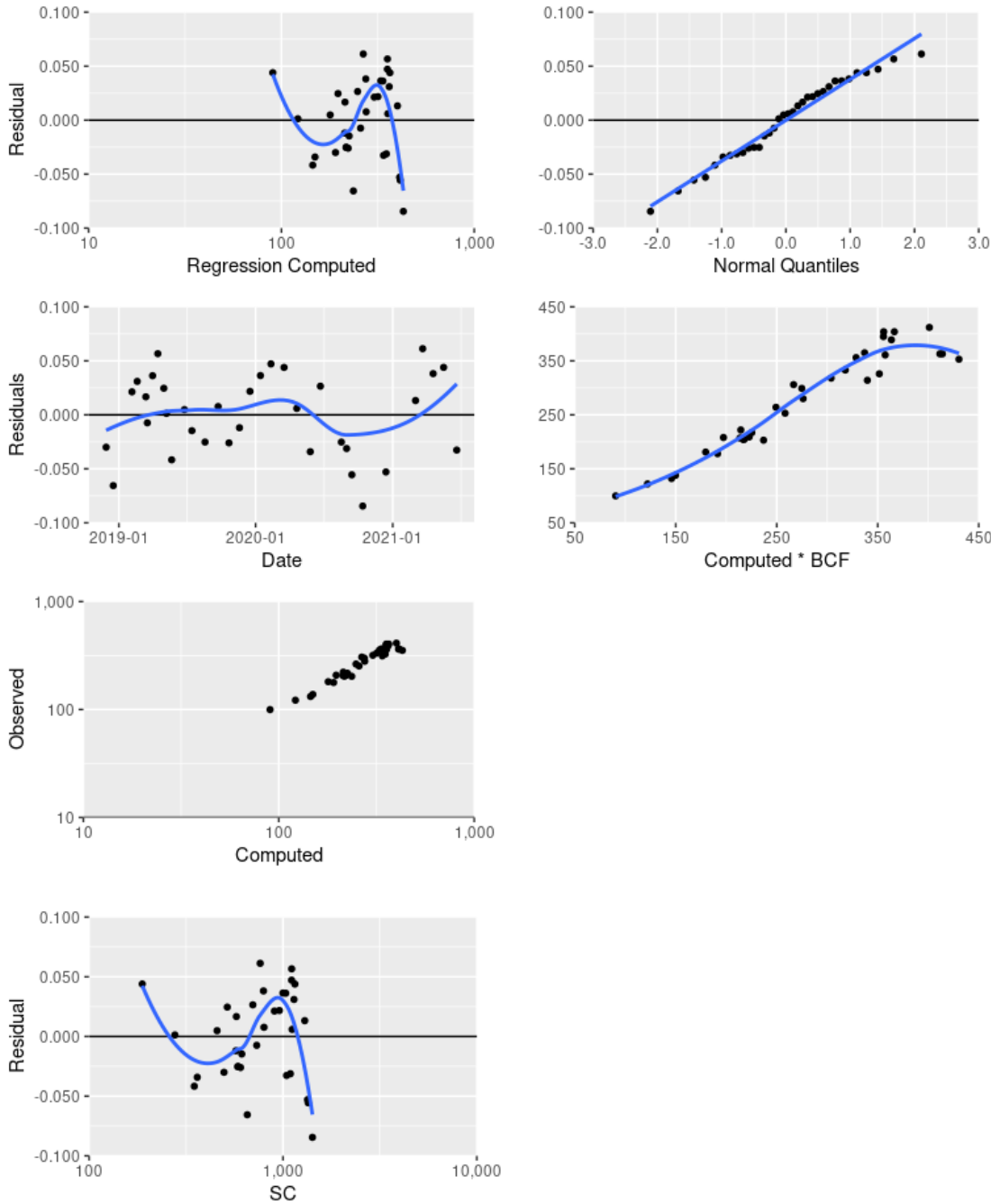
Outlier Test Criteria

Leverage	Cook's D	DFFITS
0.176	0.194	0.485

Flagged Observations

	logCaCO3	Estimate	Residual	Standard Residual	Studentized Residual	Leverage	Cook's D	DFFITS
202010130820	2.55	2.63	-0.0845	-2.32	-2.50	0.0851	0.250	-0.763
202105170800	2.00	1.96	0.0440	1.35	1.37	0.2720	0.342	0.838

Statistical Plots



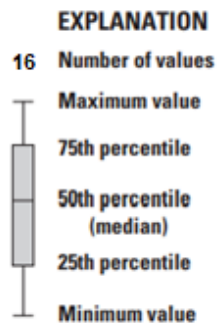
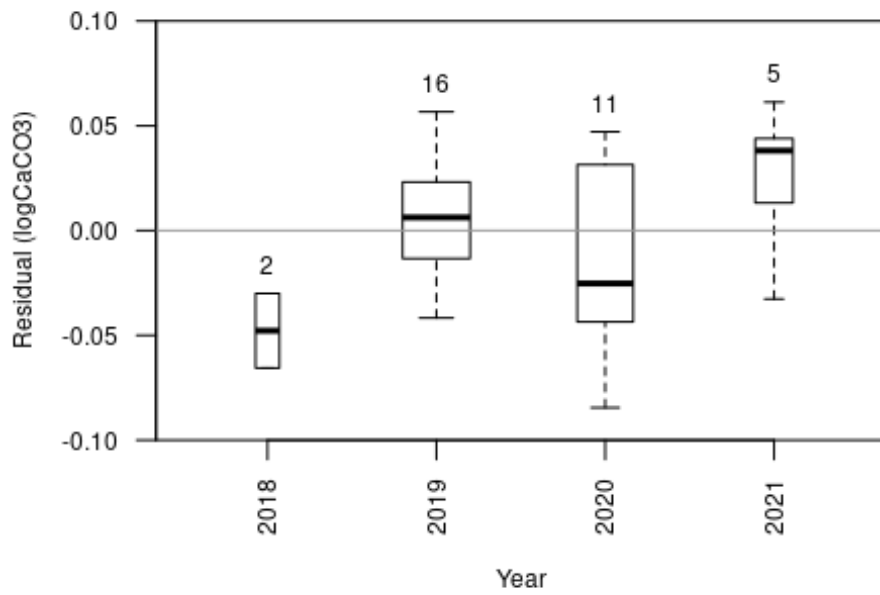
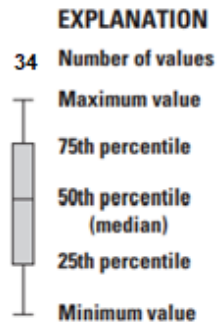
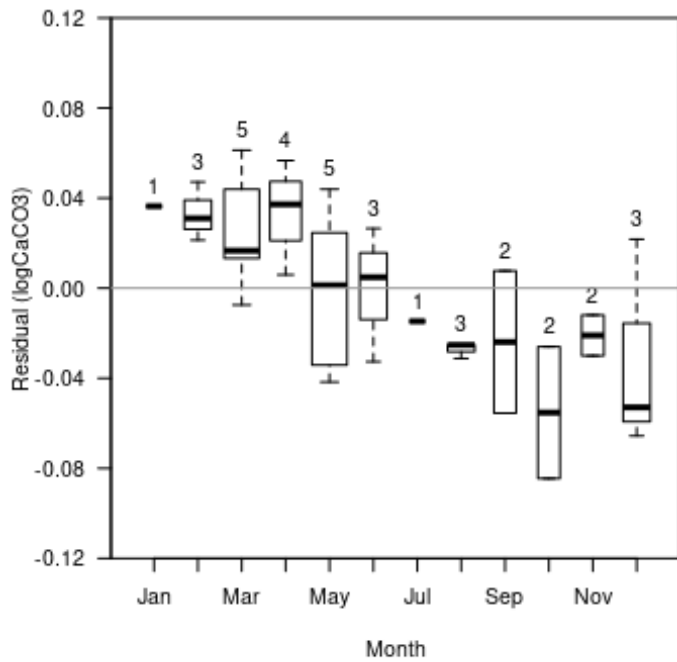
First row (left): Residual CaCO_3 related to regression computed CaCO_3 with local polynomial regression fitting, or locally estimated scatterplot smoothing (LOESS), indicated by the blue line.

First row (right): Residual CaCO_3 related to the corresponding normal quantile of the residual with simple linear regression, indicated by the blue line.

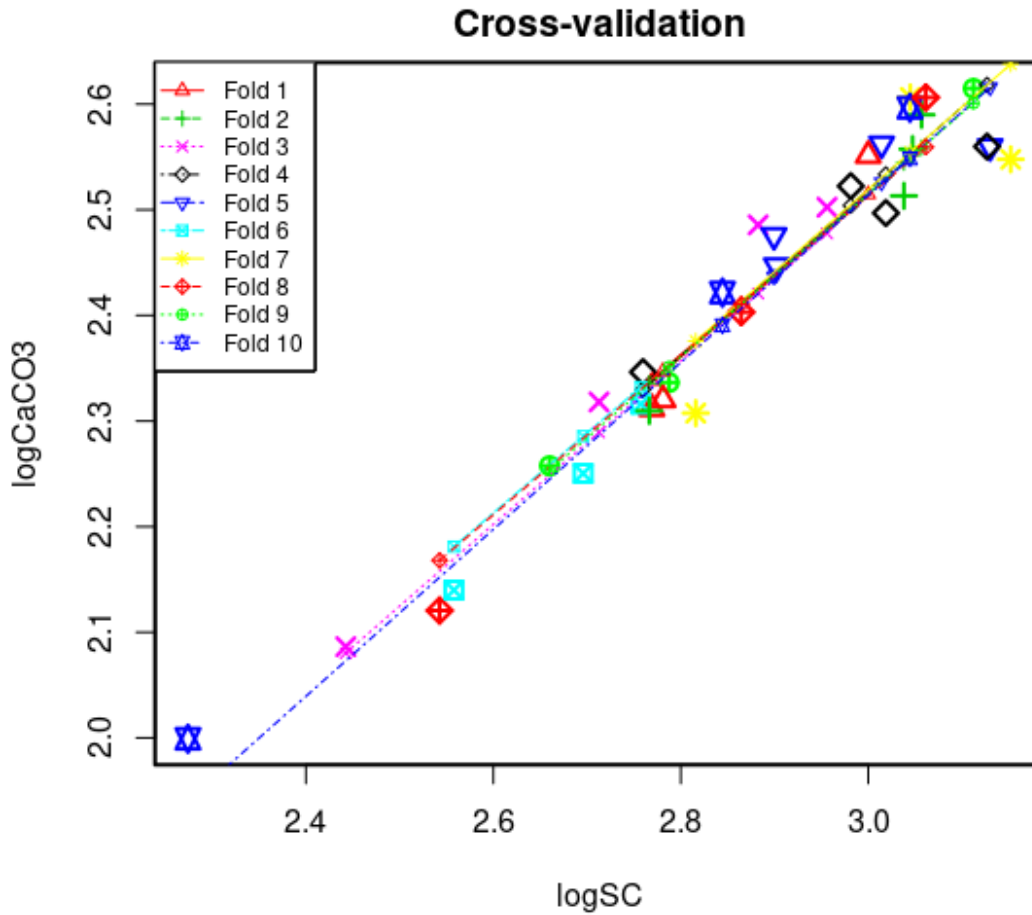
Second row: Residual CaCO_3 related to date (left) and regression computed CaCO_3 multiplied by the BCF (right) with LOESS, indicated by the blue line.

Third row: Observed CaCO_3 related to regression computed CaCO_3 .

Fourth row: Residual CaCO_3 related to SC with LOESS, indicated by the blue line.



Cross Validation



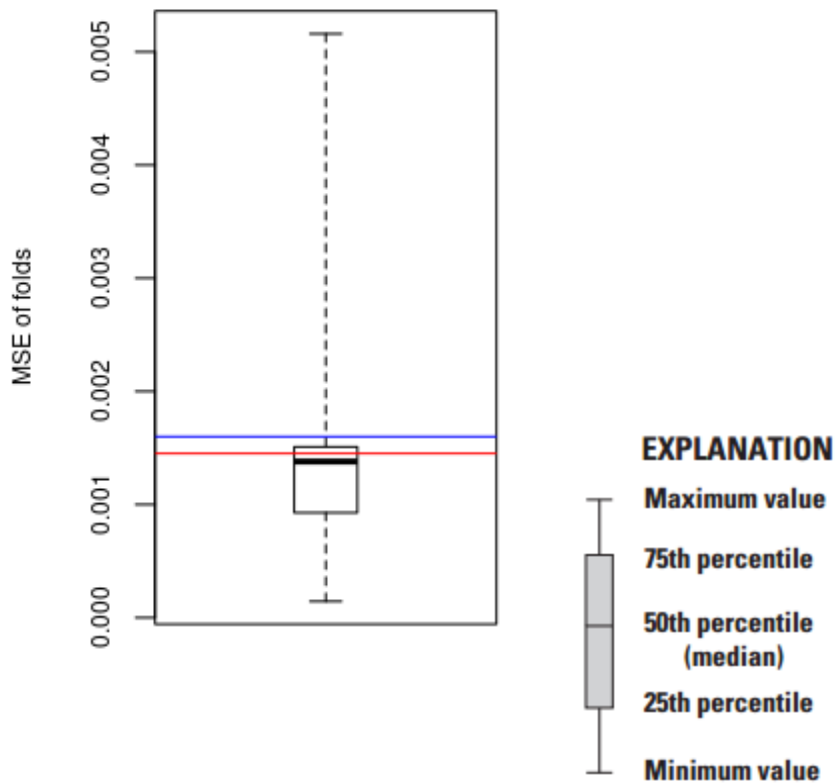
Fold - equal partition of the data (10 percent of the data).

Large symbols - observed value of a data point removed in a fold.

Small symbols - recomputed value of a data point removed in a fold.

Recomputed regression lines - adjusted regression line with one fold removed.

```
Minimum MSE of folds: 0.000144
Mean MSE of folds: 0.001600
Median MSE of folds: 0.001380
Maximum MSE of folds: 0.005160
(Mean MSE of folds) / (Model MSE): 1.100000
```



Red line - Model MSE

Blue line - Mean MSE of folds

Model-Calibration Dataset

	Date	logCaCO3	logSC	CaCO3	SC	Computed logCaCO3	Computed CaCO3	Residual	Normal Quantiles	Censored Values
0										
1	2018-11-28	2.25	2.7	178	496	2.28	191	-0.03	-0.67	--
2	2018-12-17	2.31	2.82	203	655	2.37	237	-0.0656	-1.68	--
3	2019-02-05	2.5	2.96	318	904	2.48	304	0.0213	0.336	--
4	2019-02-19	2.59	3.06	389	1140	2.56	364	0.031	0.67	--
5	2019-03-14	2.35	2.76	222	575	2.33	214	0.0167	0.259	--
6	2019-03-18	2.4	2.86	253	732	2.41	258	-0.00742	-0.184	--
7	2019-04-01	2.56	3.01	365	1030	2.53	337	0.0363	0.765	--
8	2019-04-15	2.61	3.04	404	1110	2.55	356	0.0567	1.68	--
9	2019-05-01	2.32	2.71	208	516	2.29	197	0.0246	0.496	--
10	2019-05-08	2.09	2.44	122	277	2.09	122	0.00128	-0.11	--
11	2019-05-22	2.12	2.54	132	349	2.16	146	-0.0417	-1.11	--
12	2019-06-25	2.26	2.66	181	457	2.25	180	0.00485	-0.0367	--
13	2019-07-15	2.34	2.79	217	613	2.35	225	-0.0147	-0.336	--
14	2019-08-19	2.31	2.77	205	588	2.34	218	-0.0252	-0.415	--
15	2019-09-23	2.45	2.9	280	798	2.44	276	0.0077	0.11	--
16	2019-10-22	2.32	2.78	209	604	2.35	223	-0.026	-0.581	--
17	2019-11-19	2.32	2.76	207	572	2.33	214	-0.0119	-0.259	--
18	2019-12-17	2.52	2.98	333	958	2.5	318	0.0218	0.415	--
19	2020-01-14	2.55	3	356	1000	2.52	329	0.0364	0.867	--
20	2020-02-11	2.6	3.04	395	1110	2.55	356	0.0472	1.43	--

21	2020-03-17	2.61	3.06	404	1150	2.56	366	0.0439	1.11	--
22	2020-04-20	2.56	3.05	361	1120	2.55	357	0.00596	0.0367	--
23	2020-05-26	2.14	2.56	138	361	2.17	150	-0.0342	-0.979	--
24	2020-06-22	2.42	2.84	264	699	2.4	249	0.0265	0.581	--
25	2020-08-17	2.31	2.77	204	584	2.33	217	-0.0252	-0.496	--
26	2020-08-31	2.51	3.04	326	1090	2.54	352	-0.0312	-0.765	--
27	2020-09-14	2.56	3.13	363	1350	2.62	414	-0.0555	-1.43	--
28	2020-10-13	2.55	3.15	353	1420	2.63	430	-0.0845	-2.11	--
29	2020-12-14	2.56	3.13	363	1340	2.61	412	-0.0529	-1.25	--
30	2021-03-03	2.61	3.11	412	1290	2.6	401	0.0132	0.184	--
31	2021-03-22	2.49	2.88	306	763	2.42	267	0.0613	2.11	--
32	2021-04-19	2.48	2.9	299	793	2.44	275	0.0382	0.979	--
33	2021-05-17	2	2.27	99.8	188	1.96	90.5	0.044	1.25	--
34	2021-06-21	2.5	3.02	314	1040	2.53	340	-0.0326	-0.867	--

Definitions

CaCO₃: Total hardness, in milligrams per liter as calcium carbonate (USGS parameter code 00900).

Cook's D: Cook's distance (Helsel and others, 2020).

DIFFITS: Difference in fits statistic (Helsel and others, 2020).

E.vars: Explanatory variables.

Leverage: An outlier's measure in the x direction (Helsel and others, 2020).

LOESS: Local polynomial regression fitting, or locally estimated scatterplot smoothing (Helsel and others, 2020).

LOWESS: Locally weighted scatterplot smoothing (Cleveland, 1979; Helsel and others, 2020).

MSE: Mean square error (Helsel and others, 2020).

MSPE: Model standard percentage error (Helsel and others, 2020).

Probability(>|t|): The probability that the independent variable has no effect on the dependent variable (Helsel and others, 2020).

RMSE: Root mean square error (Helsel and others, 2020).

SC: Specific conductance, in microsiemens per centimeter at 25 degrees Celsius (USGS parameter code 00095).

t value: Student's t value; the coefficient divided by its associated standard error (Helsel and others, 2020).

References Cited

Cleveland, W.S., 1979, Robust locally weighted regression and smoothing scatterplots: Journal of the American

Statistical Association, v. 74, no. 368, p. 829–836.

Cook, R.D., 1977, Detection of influential observations in linear regression: Technometrics, v. 19, no. 1, p. 15–

18. [Also available at <https://doi.org/10.2307/1268249>.]

- Duan, N., 1983, Smearing estimate—A nonparametric retransformation method: *Journal of the American Statistical Association*, v. 78, no. 383, p. 605–610. [Also available at <https://doi.org/10.1080/01621459.1983.10478017>.]
- Foster, G.M., and Graham, J.L., 2016, Logistic and linear regression model documentation for statistical relations between continuous real-time and discrete water-quality constituents in the Kansas River, Kansas, July 2012 through June 2015: U.S. Geological Survey Open-File Report 2016–1040, 27 p., accessed February 2022 at <https://doi.org/10.3133/ofr20161040>.
- Helsel, D.R., Hirsch, R.M., Ryberg, K.R., Archfield, S.A., and Gilroy, E.J., 2020, Statistical methods in water resources: U.S. Geological Survey Techniques and Methods, book 4, chap. A3, 458 p. [Also available at <https://doi.org/10.3133/tm4a3>.] [Supersedes USGS Techniques of Water-Resources Investigations, book 4, chap. A3, ver. 1.1.]
- Hem, J.D., 1985, Study and interpretation of the chemical characteristics of natural water (3d ed.): U.S. Geological Survey Water-Supply Paper 2254, 264 p. [Also available at <https://doi.org/10.3133/wsp2254>.]
- Pardoe, I., 2020, Applied regression modeling: United Kingdom, John Wiley & Sons, 336 p.
- R Core Team, 2022, R—A language and environment for statistical computing, version 4.1.2: Vienna, Austria, R Foundation for Statistical Computing, accessed December 2021 at <https://www.R-project.org/>.
- Rasmussen, P.P., Gray, J.R., Glysson, G.D., and Ziegler, A.C., 2009, Guidelines and procedures for computing time-series suspended-sediment concentrations and loads from in-stream turbidity sensor and streamflow data: U.S. Geological Survey Techniques and Methods, book 3, chap. C4, 53 p. [Also available at <https://doi.org/10.3133/tm3C4>.]
- Rasmussen, T.J., Ziegler, A.C., and Rasmussen, P.P., 2005, Estimation of constituent concentrations, densities, loads, and yields in lower Kansas River, northeast Kansas, using regression models and continuous water-quality monitoring, January 2000 through December 2003: U.S. Geological Survey Scientific Investigations Report 2005–5165, 117 p. [Also available at <https://doi.org/10.3133/sir20055165>.]
- U.S. Geological Survey, 2022, USGS water data for the Nation: U.S. Geological Survey National Water Information System database, accessed February 2022 at <https://doi.org/10.5066/F7P55KJN>.

U.S. Geological Survey, [variously dated], National field manual for the collection of water-quality data: U.S. Geological Survey Techniques of Water-Resources Investigations, book 9, chaps. A1–A9 [variously paged], accessed February 2022 at <https://water.usgs.gov/owq/FieldManual/>.

Wagner, R.J., Boulger, R.W., Jr., Oblinger, C.J., and Smith, B.A., 2006, Guidelines and standard procedures for continuous water-quality monitors—Station operation, record computation, and data reporting: U.S. Geological Survey Techniques and Methods, book 1, chap. D3, 51 p. plus 8 attachments. [Also available at <https://doi.org/10.3133/tm1D3>.] [Supersedes USGS Water-Resources Investigations Report 2000–4252.]

Williams, T.J., 2021, Linear regression model documentation and updates for computing water-quality constituent concentrations or densities using continuous real-time water-quality data for the Kansas River, Kansas, July 2012 through September 2019: U.S. Geological Survey Open-File Report 2021–1018, 18 p., accessed April 2022 at <https://doi.org/10.3133/ofr20211018>.