

# Introduction to the dataRetrieval package

Laura De Cicco<sup>1</sup> and Robert Hirsch<sup>1</sup>

<sup>1</sup>*United States Geological Survey*

January 23, 2013

## Contents

<b>1</b>	<b>Introduction to dataRetrieval</b>	<b>1</b>
<b>2</b>	<b>Getting Started</b>	<b>1</b>
2.1	New to R? . . . . .	1
2.2	R User: Installing dataRetrieval from downloaded binary . . . . .	2
2.3	R Developers: Installing dataRetrieval from gitHub . . . . .	2
<b>3</b>	<b>USGS Web Retrieval Examples</b>	<b>4</b>
3.1	USGS Web Retrieval Introduction . . . . .	4
3.2	USGS Site Information Retrievals . . . . .	5
3.3	USGS Parameter Information Retrievals . . . . .	6
3.4	USGS Daily Value Retrievals . . . . .	6
3.5	USGS Unit Value Retrievals . . . . .	9
3.6	USGS Water Quality Retrievals . . . . .	11
<b>4</b>	<b>Polished Data: USGS Web Retrieval Examples</b>	<b>14</b>

## 1 Introduction to dataRetrieval

The dataRetrieval package was created to simplify the process of getting hydrologic data in the R environment. It has been specifically designed to work seamlessly with the EGRET

package: Exploration and Graphics for RivEr Trends (EGRET). See: <https://github.com/USGS-R/EGRET/wiki> for information on EGRET.

There is a plethora of hydrological data available on the web. This package is designed specifically to load United States Geological Survey (USGS) hydrologic data to the R environment. This includes daily values, real-time (unit values), site information, and water quality sample data.

## 2 Getting Started

This section describes the options for downloading and installing the dataRetrieval package.

### 2.1 New to R?

If you are new to R, you will need to first install the latest version of R, which can be found here: <http://www.r-project.org/>.

There are many options for running and editing R code, one nice environment to learn R is RStudio. RStudio can be downloaded here: <http://rstudio.org/>. Once R and RStudio are installed, the dataRetrieval package needs to be installed as described in the next section.

### 2.2 R User: Installing dataRetrieval from downloaded binary

The latest dataRetrieval package build is available for download at [https://github.com/USGS-R/dataRetrieval/blob/master/dataRetrieval\\_1.2.1.tar.gz](https://github.com/USGS-R/dataRetrieval/blob/master/dataRetrieval_1.2.1.tar.gz). If the package's tar.gz file is saved in R's working directory, then the following command will fully install the package:

```
> install.packages("dataRetrieval_1.2.1.tar.gz", repos = NULL,  
+ type = "source")
```

If the downloaded file is stored in an alternative location, include the path in the install command. A Windows example looks like this (notice the direction of the slashes, they are in the opposite direction that Windows normally creates paths):

```
> install.packages("C:/RPackages/Statistics/dataRetrieval_1.2.1.tar.gz",  
+ repos = NULL, type = "source")
```

A Mac example looks like this:

```
> install.packages("/Users/userA/RPackages/Statistic/dataRetrieval_1.2.1.tar.gz",  
+ repos = NULL, type = "source")
```

It is a good idea to re-start the R environment after installing the package, especially if installing an updated version (that is, restart RStudio). Some users have found it necessary to delete the previous version's package folder before installing newer version of dataRetrieval. If you are experiencing issues after updating a package, trying deleting the package folder - the default location for Windows is something like this: C:/Users/userA/Documents/R/win-library/2.15/dataRetrieval, and the default for a Mac: /Users/userA/Library/R/2.15/library/dataRetrieval. Then, re-install the package using the directions above. Moving to CRAN should solve this problem.

After installing the package, you need to open the library each time you re-start R. This is done with the simple command:

```
> library(dataRetrieval)
```

Using RStudio, you could alternatively click on the checkbox for dataRetrieval in the Packages window.

### 2.3 R Developers: Installing dataRetrieval from gitHub

Alternatively, R-developers can install the latest version of dataRetrieval directly from gitHub using the devtools package. devtools is available on CRAN. Simply type the following commands into R to install the latest version of dataRetrieval available on gitHub. Rtools (for Windows) and latex tools are required.

```
> library(devtools)
> install_github("dataRetrieval", "USGS-R")
```

To then open the library, simply type:

```
> library(dataRetrieval)
```

### 3 USGS Web Retrieval Examples

In this section, we will run through 5 examples, documenting how to get raw data from the web. This includes historical daily values, real-time current values, water quality data, site information, and measured parameter information. We will use the Choptank River near Greensboro, MD as an example. The site-ID for this gage station is 01491000. Daily discharge measurements are available as far back as 1948. Additionally, forms of nitrate have been measured dating back to 1964. The functions/examples in this section are for raw data retrieval. This may or may not be the easiest data to work with. In the next section, we will use functions that retrieve and process the data in a dataframe very friendly for R analysis.

#### 3.1 USGS Web Retrieval Introduction

The United States Geological Survey organizes their hydrological data in fairly standard structure. Gage stations are located throughout the United States, each station has a unique ID. Often (but not always), these ID's are 8 digits. The first step to finding data is discovering this 8-digit ID. One potential tool for discovering data is Environmental Data Discovery and Transformation (EnDDaT): <http://cida.usgs.gov/enddat/>. Follow the example in the User's Guide to learn how to discover USGS stations and available data from any location in the United States. Essentially, you can create a Project Location on the map, set a bounding box (in miles), then search for USGS Time Series and USGS Water Quality Data. Locations, ID's, available data, and available time periods will load on the map and appropriate tabs.

Once the site-ID is known, the next required input for USGS data retrievals is the 'parameter code'. This is a 5-digit code that specifies what measured parameter is being requested. A complete list of possible USGS parameter codes can be found here:

[http://nwis.waterdata.usgs.gov/usa/nwis/pmcodes?radio\\_pm\\_search=param\\_group&pm\\_group=All+--+include+all+parameter+groups&pm\\_search=&casrn\\_search=&srsname\\_search=&format=html\\_table&show=parameter\\_group\\_nm&show=parameter\\_nm&show=casrn&show=srsname&show=parameter\\_units](http://nwis.waterdata.usgs.gov/usa/nwis/pmcodes?radio_pm_search=param_group&pm_group=All+--+include+all+parameter+groups&pm_search=&casrn_search=&srsname_search=&format=html_table&show=parameter_group_nm&show=parameter_nm&show=casrn&show=srsname&show=parameter_units)

Not every station will measure all parameters. The following is a list of commonly measured parameters:

Table 1: Commonly found USGS Parameter Codes

	pCode	shortName
1	00060	Discharge [cfs]
2	00065	Gage height [ft]
3	00010	Temperature [C]
4	00045	Precipitation [in]
5	00400	pH

For real-time data, the parameter code and site ID will suffice. The USGS stores historical

data as daily values however. The statistical process used to store the daily data is the final requirement for daily value retrievals. A 5-digit 'stat code' specifies the requested processing. A complete list of possible USGS stat codes can be found here:

[http://nwis.waterdata.usgs.gov/nwis/help/?read\\_file=stat&format=table](http://nwis.waterdata.usgs.gov/nwis/help/?read_file=stat&format=table)

The most common stat codes are:

Table 2: Commonly found USGS Stat Codes

	StatCode	shortName
1	00001	Maximum
2	00002	Minimum
3	00003	Mean
4	00008	Median

### 3.2 USGS Site Information Retrievals

To obtain all of the available site information, use the `getSiteFileData` function:

```
> siteNumber <- "01491000"
> ChopTankInfo <- getSiteFileData(siteNumber)
```

The available data for these for the USGS sites are:

```
> colnames(ChopTankInfo)

[1] "agency.cd"           "site.no"           "station.nm"
[4] "site.tp.cd"          "lat.va"            "long.va"
[7] "dec.lat.va"          "dec.long.va"       "coord.meth.cd"
[10] "coord.acy.cd"        "coord.datum.cd"    "dec.coord.datum.cd"
[13] "district.cd"         "state.cd"          "county.cd"
[16] "country.cd"          "land.net.ds"       "map.nm"
[19] "map.scale.fc"        "alt.va"            "alt.meth.cd"
[22] "alt.acy.va"          "alt.datum.cd"      "huc.cd"
[25] "basin.cd"            "topo.cd"           "instruments.cd"
[28] "construction.dt"     "inventory.dt"      "drain.area.va"
[31] "contrib.drain.area.va" "tz.cd"             "local.time.fg"
[34] "reliability.cd"      "gw.file.cd"        "nat.aqfr.cd"
[37] "aqfr.cd"             "aqfr.type.cd"      "well.depth.va"
[40] "hole.depth.va"       "depth.src.cd"      "project.no"
[43] "queryTime"
```

Pulling out a specific example piece of information, in this case station name can be done as follows:

```
> ChopTankInfo$station.nm
```

```
[1] "CHOPTANK RIVER NEAR GREENSBORO, MD"
```

Site information is obtained from <http://waterservices.usgs.gov/rest/Site-Test-Tool.html>

### 3.3 USGS Parameter Information Retrievals

To obtain all of the available information concerning a measured parameter, use the `getParameterInfo` function:

```
> parameterCd <- "00618"
> parameterINFO <- getParameterInfo(parameterCd)
```

The available data for these for the USGS sites are:

```
> colnames(parameterINFO)
```

```
[1] "parameter_cd"      "parameter_group_nm" "parameter_nm"
[4] "casrn"             "srsname"           "parameter_units"
```

Pulling out a specific example piece of information, in this case station name can be done as follows:

```
> parameterINFO$parameter_nm
```

```
[1] "Nitrate, water, filtered, milligrams per liter as nitrogen"
```

Parameter information is obtained from <http://nwis.waterdata.usgs.gov/nwis/pmcodes/>

### 3.4 USGS Daily Value Retrievals

To obtain historic daily records of USGS data, use the `retrieveNWISData` function. The arguments for this function are `siteNumber`, `parameterCd`, `startDate`, `endDate`, `statCd`, and a logical (`true/false`) `interactive`. There are 2 default argument: `statCd` defaults to "00003" and

interactive defaults to TRUE. If you want to use the default values, you do not need to list them in the function call. Setting the 'interactive' option to true will walk you through the function. It might make more sense to run large batch collections with the interactive option set to FALSE.

The dates (start and end) need to be in the format "YYYY-MM-DD". Setting the start date to "" will indicate to the program to ask for the earliest date, setting the end date to "" will ask for the latest available date.

```
> siteNumber <- "01491000"
> parameterCd <- "00060"
> startDate <- ""
> endDate <- ""
> discharge <- retrieveNWISData(siteNumber, parameterCd, startDate,
+                               endDate)
```

A dataframe is returned that looks like the following:

	agency_cd	site_no	datetime	X02_00060_00003	X02_00060_00003_cd
1	USGS	01491000	1948-01-01	190	A
2	USGS	01491000	1948-01-02	900	A
3	USGS	01491000	1948-01-03	480	A
4	USGS	01491000	1948-01-04	210	A
5	USGS	01491000	1948-01-05	210	A
6	USGS	01491000	1948-01-06	220	A

The structure of the dataframe is:

```
'data.frame':      23764 obs. of  5 variables:
 $ agency_cd      : chr  "USGS" "USGS" "USGS" "USGS" ...
 $ site_no        : chr  "01491000" "01491000" "01491000" "01491000" ...
 $ datetime       : Date, format: "1948-01-01" "1948-01-02" ...
 $ X02_00060_00003 : num  190 900 480 210 210 220 160 130 120 100 ...
 $ X02_00060_00003_cd: chr  "A" "A" "A" "A" ...
```

Note that dateTime is automatically imported as a Date. Each requested parameter has a value and remark code column. The names of these columns depend on the requested parameter and stat code combinations. USGS remark codes are often "A" (approved for publication) or "P" (provisional data subject to revision). A more complete list of remark codes can be found here: [http://waterdata.usgs.gov/usa/nwis/help?codes\\_help](http://waterdata.usgs.gov/usa/nwis/help?codes_help)

Another example that doesn't use the defaults would be a request for mean and maximum daily temperature and discharge in early 2012:

```

> siteNumber <- "01491000"
> parameterCd <- "00010,00060"
> statCd <- "00001,00003"
> startDate <- "2012-01-01"
> endDate <- "2012-06-30"
> temperatureAndFlow <- retrieveNWISData(siteNumber, parameterCd,
+   startDate, endDate, StatCd = statCd, interactive = FALSE)

```

Daily data is pulled from <http://waterservices.usgs.gov/rest/DV-Test-Tool.html>.

An example of plotting the above data (Figure 1):

```

> with(temperatureAndFlow, plot(datetime, X01_00010_00003, xlab = "Date",
+   ylab = "Temperature [C]"))
> par(new = TRUE)
> with(temperatureAndFlow, plot(datetime, X02_00060_00003, col = "red",
+   type = "l", xaxt = "n", yaxt = "n", xlab = "", ylab = "",
+   axes = FALSE))
> axis(4, col = "red", col.axis = "red")
> mtext("Discharge [cfs]", side = 4, line = 3, col = "red")

```



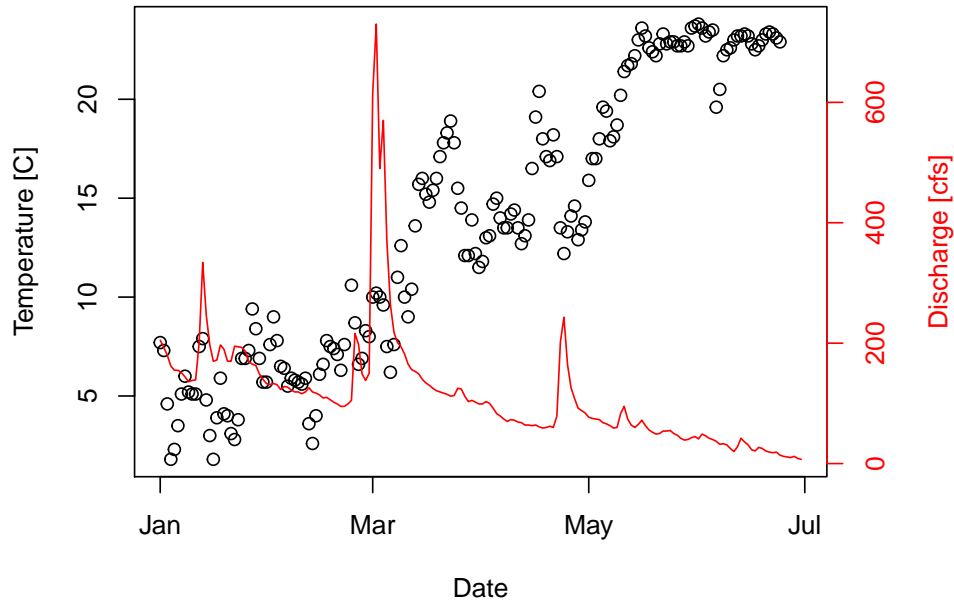


Figure 1: Temperature and discharge plot of Choptank River in 2012.

There are occasions where NWIS values are not reported as numbers, instead there might be text describing a certain event such as "Ice". Any value that cannot be converted to a number will be reported as NA in this package.

### 3.5 USGS Unit Value Retrievals

We can also get real-time, instantaneous measurements using the `retrieveUnitNWISData` function:

```
> siteNumber <- "01491000"
> parameterCd <- "00060"
> startDate <- as.character(Sys.Date() - 1)
> endDate <- as.character(Sys.Date())
> dischargeToday <- retrieveUnitNWISData(siteNumber, parameterCd,
  startDate, endDate)
```

Which produces the following dataframe:

	agency_cd	site_no	datetime	tz_cd	X02_00060	X02_00060_cd
1	USGS	01491000	2013-01-22 00:00:00	EST	209	P
2	USGS	01491000	2013-01-22 00:15:00	EST	209	P
3	USGS	01491000	2013-01-22 00:30:00	EST	209	P

4	USGS	01491000	2013-01-22 00:45:00	EST	209	P
5	USGS	01491000	2013-01-22 01:00:00	EST	209	P
6	USGS	01491000	2013-01-22 01:15:00	EST	209	P

The structure of the dataframe is:

```
'data.frame':      166 obs. of  6 variables:
 $ agency_cd   : chr  "USGS" "USGS" "USGS" "USGS" ...
 $ site_no     : chr  "01491000" "01491000" "01491000" "01491000" ...
 $ datetime    : POSIXct, format: "2013-01-22 00:00:00" "2013-01-22 00:15:00" ...
 $ tz_cd       : chr  "EST" "EST" "EST" "EST" ...
 $ X02_00060   : num  209 209 209 209 209 209 211 206 206 206 ...
 $ X02_00060_cd: chr  "P" "P" "P" "P" ...
```

Note that time now becomes important, so the `dateTime` is a `POSIXct`, and the time zone is included. Data is pulled from <http://waterservices.usgs.gov/rest/IV-Test-Tool.html>. There are occasions where NWIS values are not reported as numbers, instead a common example is "Ice". Any value that cannot be converted to a number will be reported as NA in this package.

A simple plotting example is shown in Figure 2:

```
> with(dischargeToday, plot(dateTime, X02_00060, xlab = "Date/Time",
  ylab = "Discharge [cfs]"))
```

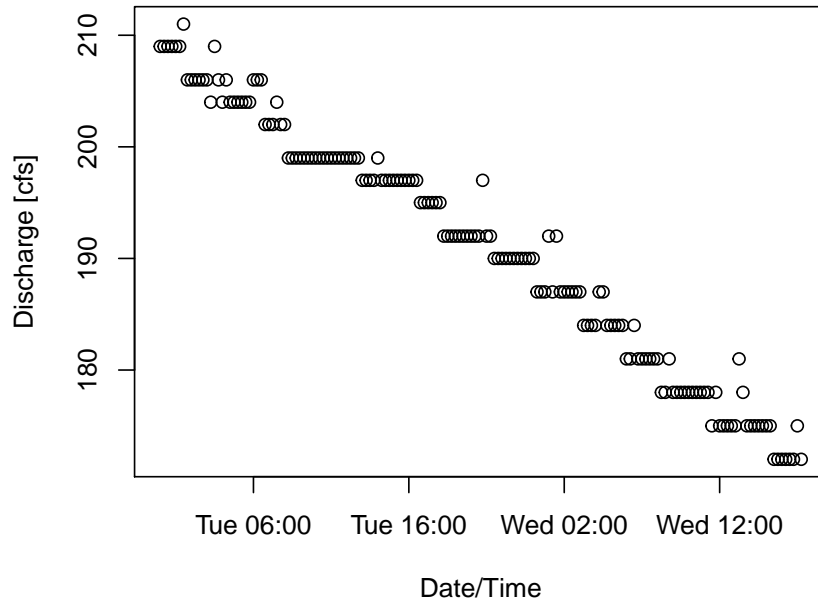


Figure 2: Real-time discharge plot of Choptank River.

### 3.6 USGS Water Quality Retrievals

Finally, we can use the `dataRetrieval` package to get water quality data that is available on the water quality data portal: <http://www.waterqualitydata.us/>. The raw data is obtained from the function `getRawQWData`, with the similar input arguments: `siteNumber`, `parameterCd`, `startDate`, `endDate`, and `interactive`. The difference is in `parameterCd`, in this function multiple parameters can be queried using a ";" separator, and setting `parameterCd <- ""` will return all of the measured observations. The raw data can be overwhelming (as will be demonstrated), a simplified version of the data can be obtained using `getQWData`.

```
> siteNumber <- "01491000"
> parameterCd <- "00618;71851"
> startDate <- "1964-06-11"
> endDate <- "2012-12-18"
> dissolvedNitrate <- getRawQWData(siteNumber, parameterCd, startDate,
  endDate)
```

Producing a dataframe with the following columns:

```
[1] "OrganizationIdentifier"
[2] "OrganizationFormalName"
[3] "ActivityIdentifier"
```

[4] "ActivityTypeCode"  
[5] "ActivityMediaName"  
[6] "ActivityMediaSubdivisionName"  
[7] "ActivityStartDate"  
[8] "ActivityStartTime.Time"  
[9] "ActivityStartTime.TimeZoneCode"  
[10] "ActivityEndDate"  
[11] "ActivityEndTime.Time"  
[12] "ActivityEndTime.TimeZoneCode"  
[13] "ActivityDepthHeightMeasure.MeasureValue"  
[14] "ActivityDepthHeightMeasure.MeasureUnitCode"  
[15] "ActivityDepthAltitudeReferencePointText"  
[16] "ActivityTopDepthHeightMeasure.MeasureValue"  
[17] "ActivityTopDepthHeightMeasure.MeasureUnitCode"  
[18] "ActivityBottomDepthHeightMeasure.MeasureValue"  
[19] "ActivityBottomDepthHeightMeasure.MeasureUnitCode"  
[20] "ProjectIdentifier"  
[21] "ActivityConductingOrganizationText"  
[22] "MonitoringLocationIdentifier"  
[23] "ActivityCommentText"  
[24] "SampleAquifer"  
[25] "HydrologicCondition"  
[26] "HydrologicEvent"  
[27] "SampleCollectionMethod.MethodIdentifier"  
[28] "SampleCollectionMethod.MethodIdentifierContext"  
[29] "SampleCollectionMethod.MethodName"  
[30] "SampleCollectionEquipmentName"  
[31] "ResultDetectionConditionText"  
[32] "CharacteristicName"  
[33] "ResultSampleFractionText"  
[34] "ResultMeasureValue"  
[35] "ResultMeasure.MeasureUnitCode"  
[36] "MeasureQualifierCode"  
[37] "ResultStatusIdentifier"  
[38] "StatisticalBaseCode"  
[39] "ResultValueTypeName"  
[40] "ResultWeightBasisText"  
[41] "ResultTimeBasisText"  
[42] "ResultTemperatureBasisText"  
[43] "ResultParticleSizeBasisText"  
[44] "PrecisionValue"  
[45] "ResultCommentText"  
[46] "USGSPCode"  
[47] "ResultDepthHeightMeasure.MeasureValue"  
[48] "ResultDepthHeightMeasure.MeasureUnitCode"

```

[49] "ResultDepthAltitudeReferencePointText"
[50] "SubjectTaxonomicName"
[51] "SampleTissueAnatomyName"
[52] "ResultAnalyticalMethod.MethodIdentifier"
[53] "ResultAnalyticalMethod.MethodIdentifierContext"
[54] "ResultAnalyticalMethod.MethodName"
[55] "MethodDescriptionText"
[56] "LaboratoryName"
[57] "AnalysisStartDate"
[58] "ResultLaboratoryCommentText"
[59] "DetectionQuantitationLimitTypeName"
[60] "DetectionQuantitationLimitMeasure.MeasureValue"
[61] "DetectionQuantitationLimitMeasure.MeasureUnitCode"
[62] "PreparationStartDate"

```

To get a simplified dataframe that contains only datetime, value, and qualifier, use the function `getQWData`:

```

> dissolvedNitrateSimple <- getQWData(siteNumber, parameterCd,
  startDate, endDate)
> head(dissolvedNitrateSimple)

```

	dateTime	qualifier.71851	value.71851	qualifier.00618	value.00618
1	1964-06-11		3.3		0.745
3	1964-09-10		5.3		1.200
5	1965-02-01		2.9		0.655
7	1965-02-25		2.4		0.542
9	1965-03-25		1.5		0.339
11	1965-04-20		2.2		0.497

Note that in this dataframe, datetime is only imported as Dates (no times are included), and the qualifier is either blank or "<" signifying a censored value.

An example of plotting the above data (Figure 3):

```

> with(dissolvedNitrateSimple, plot(dateTime, value.00618, xlab = "Date",
  ylab = paste(parameterINFO$srsname, "[", parameterINFO$parameter_units,
  "]")))

```

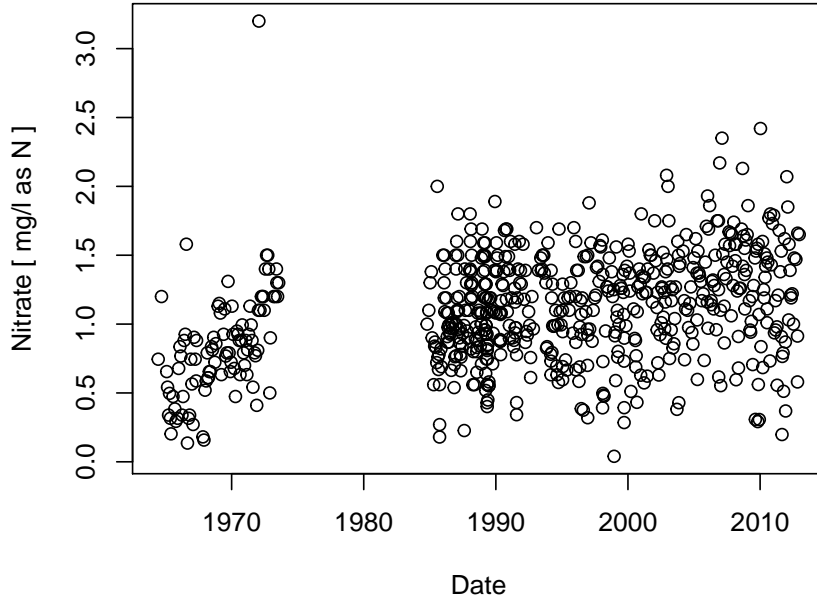


Figure 3: Nitrate plot of Choptank River.

## 4 Polished Data: USGS Web Retrieval Examples

Rather than using raw data as retrieved by the web, the `dataRetrieval` package also includes functions that return the data in a structure that has been designed to work with the EGRET R package. In general, these dataframes may be much more 'R-friendly' than the raw data, and will contain additional information that allows for efficient data analysis.

In this section, we use 3 `dataRetrieval` functions to get sufficient data to perform an EGRET analysis. We will continue analyzing the Choptank River. We will need essentially the same data that was retrieved in the previous section, but we will get the daily discharge values in a dataframe called `Daily`, the nitrate sample data in a dataframe called `Sample`, and the data about the station and parameters in a dataframe called `INFO`. These are the dataframes that were exclusively designed to work with the EGRET R package, however can be very useful for all hydrologic studies.

The function to obtain the daily values (discharge in this case) is `getDVDData`. It requires the inputs `siteNumber`, `ParameterCd`, `StartDate`, `EndDate`, `interactive`, and `convert`. Most of these arguments are described in the previous section, however 'convert' is a new argument, it's default is `TRUE`, and it tells the program to convert the values from cfs to cms. If you don't want this conversion, set `convert=FALSE` in the function call.

The function to obtain sample data from the water quality portal is `getSampleData`. The arguments for this function are also `siteNumber`, `ParameterCd`, `StartDate`, `EndDate`, `interactive`.

```

> siteNumber <- "01491000"
> parameterCd <- "00631"
> startDate <- "1964-01-01"
> endDate <- "2013-01-01"
> Daily <- getDVDData(siteNumber, "00060", startDate, endDate)

```

There are 17899 data points, and 17899 days.

There are 0 zero flow days

If there are any zero discharge days, all days had 0 cubic meters per second added to the c

```

> Sample <- getSampleData(siteNumber, parameterCd, startDate, endDate)
> INFO <- getMetaData(siteNumber, parameterCd, interactive = FALSE)
> Sample <- mergeReport()

```

Discharge Record is 17899 days long, which is 49 years

First day of the discharge record is 1964-01-01 and last day is 2013-01-01

The water quality record has 627 samples

The first sample is from 1973-06-04 and the last sample is from 2012-12-18

Discharge: Minimum, mean and maximum 0.00991 4.02 246

Concentration: Minimum, mean and maximum 0.05 1.1 2.4

Percentage of the sample values that are censored is 0.16 %

```

> head(Sample)

```

	Date	ConcLow	ConcHigh	Uncen	ConcAve	Julian	Month	Day	DecYear	MonthSeq
1	1973-06-04	1.30	1.30	1	1.30	45079	6	155	1973.422	1482
2	1979-09-25	0.52	0.52	1	0.52	47383	9	268	1979.731	1557
3	1979-10-24	0.62	0.62	1	0.62	47412	10	297	1979.810	1558
4	1979-12-05	1.40	1.40	1	1.40	47454	12	339	1979.925	1560
5	1979-12-21	1.20	1.20	1	1.20	47470	12	355	1979.969	1560
6	1980-01-24	0.84	0.84	1	0.84	47504	1	24	1980.064	1561
	SinDY	CosDY	Q	LogQ						
1	0.4699767	-0.8826788	3.256437	1.180634						
2	-0.9927882	-0.1198812	3.398022	1.223193						
3	-0.9295235	0.3687629	3.199804	1.163089						
4	-0.4547551	0.8906165	2.973269	1.089662						
5	-0.1961425	0.9805754	2.944952	1.080093						
6	0.3925740	0.9197204	10.901986	2.388945						

## References

- [1] Helsel, D.R. and R. M. Hirsch, 2002. Statistical Methods in Water Resources Techniques of Water Resources Investigations, Book 4, chapter A3. U.S. Geological Survey. 522 pages. <http://pubs.usgs.gov/twri/twri4a3/>
- [2] Hirsch, R. M., Moyer, D. L. and Archfield, S. A. (2010), Weighted Regressions on Time, Discharge, and Season (WRTDS), with an Application to Chesapeake Bay River Inputs. JAWRA Journal of the American Water Resources Association, 46: 857-880. doi: 10.1111/j.1752-1688.2010.00482.x <http://onlinelibrary.wiley.com/doi/10.1111/j.1752-1688.2010.00482.x/full>
- [3] Sprague, L. A., Hirsch, R. M., and Aulenbach, B. T. (2011), Nitrate in the Mississippi River and Its Tributaries, 1980 to 2008: Are We Making Progress? Environmental Science & Technology, 45 (17): 7209-7216. doi: 10.1021/es201221s <http://pubs.acs.org/doi/abs/10.1021/es201221s>