

Package ‘webreadr’

August 29, 2016

Type Package

Title Tools for Reading Formatted Access Log Files

Version 0.4.0

Date 2016-01-23

Author Oliver Keyes

Maintainer Oliver Keyes <ironholds@gmail.com>

Description R is used by a vast array of people for a vast array of purposes - including web analytics. This package contains functions for consuming and munging various common forms of request log, including the Common and Combined Web Log formats and various Amazon access logs.

License MIT + file LICENSE

BugReports <https://github.com/Ironholds/webreadr/issues>

URL <https://github.com/Ironholds/webreadr>

Suggests iptools, urltools, rgeolocate, knitr, testthat

LinkingTo Rcpp

Imports Rcpp, readr

VignetteBuilder knitr

RoxygenNote 5.0.1

NeedsCompilation yes

Repository CRAN

Date/Publication 2016-01-23 23:19:32

R topics documented:

read_aws	2
read_clf	3
read_combined	4
read_s3	5
read_squid	7

split_clf	8
split_squid	9
webreadr	9

Index	10
--------------	-----------

read_aws	<i>read Amazon CloudFront access logs</i>
----------	---

Description

Amazon CloudFront uses access logs with a standard format described [on their website](#). `read_aws` reads these files in; due to the Amazon treatment of header lines, it is capable of organically detecting whether files lack common fields, and compensating for that. See "Details"

Usage

```
read_aws(file)
```

Arguments

`file` the full path to the AWS file you want to read.

Details

Amazon CloudFront uses tab-separated files with [Amazon-specific fields](#). This can be changed by individual CloudFront users, however, to exclude particular fields, and historically has contained fewer fields than it now does. Luckily, Amazon's insistence on standardisation in field names means that we can organically detect if fields are missing, and compensate for that before reading in the file.

If no fields are missing, the fields returned will be:

- `date`: the date and time when the request was *completed*
- `time_elapsed`: the amount of time (in milliseconds) that the connection and fulfilment of the request lasted for.
- `edge_location`: the Amazon edge location that served the request, identified by a three-letter code. See the Amazon documentation for more details.
- `bytes_sent`: a count of the number of bytes sent by the server to the client, including headers, to fulfil the request.
- `ip_address`: the IP address of the client making the request.
- `http_method`: the HTTP method (POST, GET, etc) used.
- `host`: the CloudFront host name.
- `path`: the path to the requested asset.
- `status_code`: the HTTP status code associated with the request.
- `referer`: the referer associated with the request.

- `user_agent`: the user agent of the client that made the request.
- `query`: the query string associated with the request; if there is no query string, this will be a dash.
- `cookie`: the cookie header from the request, stored as name-value pairs. When no cookie header is provided, or it is empty, this will be a dash.
- `result_type`: the result of the request. This is similar to Squid response codes (see [read_squid](#)) but Amazon-specific; their documentation contains details on what each code means.
- `request_id`: A hashed unique identifier for each request.
- `host_header`: the host header of the requested asset. While `host` will always be the CloudFront host name, `host_header` contains alternate domain names (or 'CNAMES') when the CloudFront distribution is using them.
- `protocol`: the protocol used in the request (http/https).
- `bytes_received`: client-to-server bytes, including headers.
- `time_elapsed`: the time elapsed, in seconds, between the time the request was received and the time the server completed responding to it.

See Also

[read_s3](#), for Amazon S3 files, [read_clf](#) for the Common Log Format, [read_squid](#) and [read_combined](#).

Examples

```
#Read in an example CloudFront file provided with the webreadr package.  
data <- read_aws(system.file("extdata/log.aws", package = "webreadr"))
```

<code>read_clf</code>	<i>read CLF-formatted logs</i>
-----------------------	--------------------------------

Description

Read a file of request logs stored in the **Common Log Format**.

Usage

```
read_clf(file, has_header = FALSE)
```

Arguments

<code>file</code>	the full path to the CLF-formatted file you want to read.
<code>has_header</code>	whether or not the file has a header row. Set to FALSE by default.

Details

the CLF is a standardised format for web request logs. It consists of the fields:

- `ip_address`: the IP address of the remote host that made the request. The CLF does not (by default) include the de-facto standard X-Forwarded-For header
- `remote_user_ident`: the [RFC 1413](#) remote user identifier.
- `local_user_ident`: the identifier the user has authenticated with locally.
- `timestamp`: the timestamp associated with the request, stored as "[08/Apr/2001:17:39:04 -0800]", where "-0800" represents the time offset (minus eight hours) of the timestamp from UTC.
- `request`: the actual user request, containing the HTTP method used, the asset requested, and the HTTP Protocol version used.
- `status_code`: the HTTP status code returned.
- `bytes_sent`: the number of bytes sent

While outdated as a standard, systems using the CLF are still around; the Squid caching system, for example, uses the CLF as one of its default log formats (the other, the squid "native" format, can be read with [read_squid](#)).

Value

a data.frame consisting of seven fields, as discussed above, with normalised timestamps.

See Also

[read_combined](#) for the /Combined/ Log Format, and [split_clf](#) for splitting out the "requests" field.

Examples

```
#Read in an example CLF-formatted file provided with the webreadr package.
data <- read_clf(system.file("extdata/log.clf", package = "webreadr"))
```

<code>read_combined</code>	<i>read Combined Log Format files</i>
----------------------------	---------------------------------------

Description

read requests logs following the Combined Log Format.

Usage

```
read_combined(file, has_header = FALSE)
```

Arguments

file	the full path to the CLF-formatted file you want to read.
has_header	whether or not the file has a header row. Set to FALSE by default.

Details

the Combined Log Format (CLF) is the same as the Common Log Format (CLF, because software engineers and naming go together like chalk and cheese), which is documented at [read_clf](#). In addition to the fields described there, the Combined Log Format also includes:

- referer: the referer associated with the request.
- user_agent: the user agent of the user that made the request.

read_combined handles these fields, as well as the CLF-standard ones. This is (amongst other things) the default logging format for [nginx](#) servers

See Also

[read_clf](#) for the /Common/ Log Format, and [split_clf](#) for splitting out the "requests" field.

Examples

```
#Read in an example Combined-formatted file provided with the webreadr package.  
data <- read_combined(system.file("extdata/combined_log.clf", package = "webreadr"))
```

read_s3	<i>Read Amazon S3 Access Logs</i>
---------	-----------------------------------

Description

read_s3 provides a reader for Amazon's S3 service's access logs, described [here](#).

Usage

```
read_s3(file)
```

Arguments

file	the full path to the S3 file you want to read.
------	--

Details

S3 access logs contain information about requests to S3 buckets, and follow a standard format described [here](#).

The fields for S3 files are:

- owner: the owner of the S3 bucket; a hashed user ID
- bucket: the bucket that processed the request.
- request_time: the time that a request was received. Formatted as POSIXct timestamps.
- remote_ip: the IP address that made the request.
- requester: the user ID of the person making the request; Anonymous if the request was not authenticated.
- operation: the actual operation performed with the request.
- key: the request's key, normally an encoded URL fragment or NA if the operation did not contain a key.
- uri: the full URI for the request, as well as the HTTP method and version. `split_clf` works to split this into a `data.frame` of 3 columns.
- status: the HTTP status code associated with the request.
- error: the error code, if an error occurred; NA otherwise. See [here](#) for more information about S3 error codes.
- sent: the number of bytes returned in response to the request.
- size: the total size of the returned object.
- time: the number of milliseconds between the request being sent and the response being sent, from the server's perspective.
- turn_around: the number of milliseconds the S3 bucket spent processing the request.
- referer: the referer associated with the request.
- user_agent: the user agent associated with the request.
- version_id: the version ID of the request; NA if the requested operation does not involve a version ID.

See Also

[read_aws](#) for reading Amazon Web Services (AWS) access log files, and [split_clf](#), which works well on the `uri` field from S3 files.

Examples

```
# Using the inbuilt testing dataset
s3_data <- read_s3(system.file("extdata/s3.log", package = "webreadr"))
```

read_squid	<i>read Squid files</i>
------------	-------------------------

Description

the Squid default log formats are either the CLF - for which, use [read_clf](#) - or the "native" Squid format, which is described in more detail below. `read_squid` allows you to read the latter.

Usage

```
read_squid(file, has_header = FALSE)
```

Arguments

<code>file</code>	the full path to the CLF-formatted file you want to read.
<code>has_header</code>	whether or not the file has a header row. Set to <code>FALSE</code> by default.

Details

The log format for Squid servers can be custom-set, but by default follows one of two patterns; it's either the Common Log Format (CLF), which you can read in with [read_clf](#), or the "native log format", a Squid-specific format handled by this function. It consists of the fields:

- `timestamp`: the timestamp identifying when the request was received. This is stored (from the file's point of view) as a count of seconds, in UNIX time: `read_squid` turns them into POSIXlt timestamps, assuming UTC as an origin timezone.
- `time_elapsed`: the amount of time (in milliseconds) that the connection and fulfilment of the request lasted for.
- `ip_address`: the IP address of the remote host making the request.
- `status_code`: the status code and Squid response code associated with that request, stored as a single field. This can be split into two distinct fields with [split_squid](#)
- `bytes_sent`: the number of bytes sent
- `http_method`: the HTTP method (POST, GET, etc) used.
- `url`: the URL of the requested asset.
- `remote_user_ident`: the [RFC 1413](#) remote user identifier.
- `peer_info`: the status of how forwarding to a peer server was handled and, if the request was forwarded, the server it was sent to.

See Also

[read_clf](#) for the Common Log Format (also used by Squids), and [split_squid](#) for splitting the "status_code" field into its component parts.

Examples

```
#Read in an example Squid file provided with the webreadr package.  
data <- read_squid(system.file("extdata/log.squid", package = "webreadr"))
```

split_clf	<i>split requests from a CLF-formatted file</i>
-----------	---

Description

CLF (Combined/Common Log Format) files store the HTTP method, protocol and asset requested in the same field. `split_clf` takes this field as a vector and returns a data.frame containing these elements in distinct columns. The function also works nicely with the `uri` field from Amazon S3 files (see [read_s3](#)).

Usage

```
split_clf(requests)
```

Arguments

`requests` the "request" field from a CLF-formatted file, read in with [read_clf](#) or [read_combined](#).

Value

a data.frame of three columns - "method", "asset" and "protocol" - representing, respectively, the HTTP method used ("GET"), the asset requested ("/favicon.ico") and the protocol used ("HTTP/1.0"). In cases where the request is not intact (containing, for example, just the protocol or just the asset) a row of empty strings will currently be returned. In the future, this will be somewhat improved.

See Also

[read_clf](#) and [read_combined](#) for reading in these files.

Examples

```
# Grab CLF data and split out the request.
data <- read_combined(system.file("extdata/combined_log.clf", package = "webreadr"))
requests <- split_clf(data$request)

# An example using S3 files
s3_data <- read_s3(system.file("extdata/s3.log", package = "webreadr"))
s3_requests <- split_clf(s3_data$uri)
```

split_squid	<i>split the "status_code" field in a Squid-formatted dataset.</i>
-------------	--

Description

the Squid data format (which can be read in with [read_squid](#)) stores the squid response and the HTTP status code as a single field. [split_squid](#) allows you to split these into a data.frame of two distinct columns.

Usage

```
split_squid(status_codes)
```

Arguments

status_codes a status_code column from a Squid file read in with [read_squid](#)

Value

a data.frame of two columns - "squid_code" and "http_status" - representing, respectively, the Squid response to the request and the HTTP status of it. In cases where the status code is not intact (containing, for example, just the squid_code) a row of empty strings will currently be returned. In the future, this will be somewhat improved.

See Also

[read_squid](#) for reading these files in, and [split_clf](#) for similar parsing of multi-field columns in Common/Combined Log Format (CLF) data.

Examples

```
#Read in an example Squid file provided with the webtools package, then split out the codes
data <- read_squid(system.file("extdata/log.squid", package = "webreadr"))
statuses <- split_squid(data$status_code)
```

webreadr	<i>A package for reading various common forms of request log</i>
----------	--

Description

see the [introductory vignette](#) for more details!

Index

`read_aws`, 2, 6
`read_clf`, 3, 3, 5, 7, 8
`read_combined`, 3, 4, 4, 8
`read_s3`, 3, 5, 8
`read_squid`, 3, 4, 7, 9

`split_clf`, 4–6, 8, 9
`split_squid`, 7, 9, 9

`webreadr`, 9