

Package ‘epiforecast.cpp14funs’

September 6, 2019

Type Package

Title Tools for forecasting semi-regular seasonal epidemic curves and similar time series

Version 0.0.1

Date 2015-12-09

Author Logan C. Brooks, David C. Farrow, Sangwon Hyun, Ryan J. Tibshirani, Roni Rosenfeld

Maintainer Logan C. Brooks <lcbrooks@cs.cmu.edu>

Recommends epiforecast

Imports stats,
utils,
lubridate,
splines,
glmgen,
glmnet,
httr,
Rcpp,
rlist,
Hmisc,
weights,
hexbin,
plotrix,
tibble,
dplyr,
pipeR,
matrixStats,
quantreg,
lpSolve,
parallel,
R.utils,
RCurl,
pbmcapply

Description Tools for forecasting semi-regular seasonal epidemic curves and similar time series. Includes an empirical Bayes approach that forms a prior

by transforming historical curves, a basis regression approach that balances matching observations from the current season and matching historical seasons' measurements for future weeks, and timestep-by-timestep weighted kernel density estimation on backward differences parameterized by both the time series measurements and the current time.

License GPL-2

RoxygenNote 6.1.1

Collate 'RcppExports.R'
'eb_dists.R'
'interface.R'
'eb.R'

LinkingTo Rcpp

SystemRequirements C++14

Suggests testthat

R topics documented:

eb.createForecasts	2
eb.fitSmoothCurves	3
eb.sim	4
fit.eb.control.list	5
get_eb_control_list	6
match.dist	9
match.single.integer	10
unifChoicePrior	10

Index	12
--------------	-----------

eb.createForecasts	<i>Function for making forecasts with the empirical Bayes method.</i>
--------------------	---

Description

Function for making forecasts with the empirical Bayes method.

Usage

```
eb.createForecasts(dat, new.dat, fit.obj = NULL, time.of.forecast = NULL,  
  control.list = get.eb.control.list())
```

Arguments

<code>dat</code>	a list of numeric vectors, one per past season, containing historical trajectories.
<code>new.dat</code>	a single numeric vector containing the observations for the current season so far, and possibly future data points as well (when performing retrospective analysis); should not contain any NA's.
<code>fit.obj</code>	a collection of fit curves and noise level estimates to use when forming the prior; defaults to <code>smooth.curves.to.fit(eb.fitSmoothCurves(dat))</code> ; while the smoothing method is quite fast, repeated calls to <code>eb.createForecasts</code> may benefit from caching the smoothed curves and feeding them in each time.
<code>time.of.forecast</code>	integer in <code>[0..length(new.dat.partial)]</code> ; if specified, the forecast is prepared as if <code>new.dat.partial[seq_len(time.of.forecast)]</code> was fed in.
<code>control.list</code>	optional control list to forward to eb.createForecasts .

Value

a list with two components:

`ys`: a numeric matrix; in most other methods, each column is a different possible trajectory for the current season, with NA's in `new.dat` filled in with random draws from the forecasted distribution, and non-NA's (observed data) filled in with an imagined resampling of noise based on the model.

`weights`: a numeric vector; assigns a weight to each column of `ys`, which is used by methods relying on importance sampling.

<code>eb.fitSmoothCurves</code>	<i>Function for fitting smooth curves to past seasons' data</i>
---------------------------------	---

Description

Arguments:

Usage

```
eb.fitSmoothCurves(dat.obj, method = c("ss", "tf"), cv.rule = c("min",
  "lse", "gcv"), tf.ord = 2, verbose = FALSE)
```

Arguments

<code>dat.obj</code>	assumed to be a list, of length equal to number of past seasons. Each item here is itself a list, each component containing a vector of "signals" for that seasons
<code>method</code>	one of "ss" and "tf". The former uses R's built-in smoothing spline method; the latter uses the <code>glmgen</code> package
<code>cv.rule</code>	one of "min", "lse", or "gcv": the rule for selecting the smoothing parameter, where "min" gives the CV usual rule, and "lse" uses the CV 1-standard-error rule, and "gcv" uses generalized cross-validation (more efficient, less accurate); the "ss" method accepts only "min" and "gcv"

tf.ord	the order of the piecewise polynomial fit by trend filtering. Default is 2
verbose	logical; if TRUE, progress information will be printed out to the terminal.
cv	if TRUE, uses cross-validation to find the smoothing parameter; if FALSE, uses generalized cross-validation (more efficient, less accurate)

Value

a list with components

smooth.obj: a list of the same dimension as dat.obj, except all observed signal values have all been replaced by smoothed values

sigma.hat: a vector of length equal to the number of seasons, each component giving an estimate of the standard deviation of the noise in that season's data

eb.sim

Function for making forecasts with the empirical Bayes method.

Description

Function for making forecasts with the empirical Bayes method.

Usage

```
eb.sim(full.dat, baseline = 0, time.of.forecast = NULL,
       max.n.sims = 1000L, fit.obj = NULL,
       control.list = get_eb_control_list())
```

Arguments

max.n.sims	single non-NA integer value or NULL: the number of curves to sample from the inferred distribution
fit.obj	optional argument to forward to eb.createForecasts
control.list	optional argument to forward to eb.createForecasts ; n.out is overridden with length(new.dat)
max.n.sims	maximum number of (weighted) simulated curves to produce. Defaults to 1000.

Value

a sim object — a list with two components:

ys: a numeric matrix, typically with multiple columns; each column is a different possible trajectory for the current season, with NA's in the input for the current season filled in with random draws from the forecasted distribution, and non-NA's (observed data) filled in with an imagined resampling of noise based on the model (for some models, the non-NA values will remain unchanged).

weights: a numeric vector; assigns a weight to each column of ys, which is used by methods relying on importance sampling.

Author(s)

Logan C. Brooks, David C. Farrow, Sangwon Hyun, Ryan J. Tibshirani, Roni Rosenfeld

Examples

```
library(epiforecast)
fluvview.nat.recent.df =
  trimPartialPastSeasons(fetchEpidataDF("fluvview", "nat",
    first.week.of.season=21L,
    cache.file.prefix="fluvview_nat_allfetch"),
    "wili", min.points.in.season=52L)
## Recent historical seasons + current season, minus 2009 (nonseasonal
## pandemic) season:
full.dat = split(fluview.nat.recent.df$wili, fluview.nat.recent.df$season)
names(full.dat) <- sprintf("S%s", names(full.dat))
full.dat <- full.dat[names(full.dat)!="S2009"]
## Recent historical seasons minus 2009:
dat = head(full.dat, -1L)
## Current season:
new.dat = tail(full.dat, 1L)[[1]]
## Sample from conditional curve distribution estimate using CDC's 2015
## national %wILI onset threshold baseline of 2.1:
sim = eb.sim(dat, new.dat, 2.1, max.n.sims=100)

print(sim)
plot(sim)
```

fit.eb.control.list	<i>Takes an EB control list containing arguments that may require fitting to a curve object, performs any fitting, and outputs a static EB control list containing the results of the fitting procedure. The contents of the EB control list are also validated and standardized to a more rigid form, e.g, replacing some NULL values with defaults and some non-integer-class integral input with integer-class versions.</i>
---------------------	---

Description

Takes an EB control list containing arguments that may require fitting to a curve object, performs any fitting, and outputs a static EB control list containing the results of the fitting procedure. The contents of the EB control list are also validated and standardized to a more rigid form, e.g, replacing some NULL values with defaults and some non-integer-class integral input with integer-class versions.

Usage

```
fit.eb.control.list(dat, new.dat, fit.obj, time.of.forecast, control.list)
```

get_eb_control_list *Generates a control list for [eb.createForecasts](#).*

Description

With no arguments, returns the default control list. Optional arguments provided to the function will override these default values (currently with no validation checks).

Usage

```
get_eb_control_list(parent = NULL, max.n.sims = 2000L,
  peak.time.dist = NULL, x.shift.dist = NULL, x.scale.dist = NULL,
  y.scale.baseline = 0, peak.height.dist = NULL, y.scale.dist = NULL,
  sd.option = c("match", "scale", "prior"), sd.prior = "uniform",
  sd.scale.dist = NULL, reasonable.future.weight = 0,
  n.future.neighbors = 3L, bias_peaktime_weight = NA_real_,
  bias_peaktime_mean = NA_real_, bias_peaktime_sd = NA_real_,
  bias_peakheight_weight = NA_real_, bias_peakheight_mean = NA_real_,
  bias_peakheight_sd = NA_real_, inactive.seasons = NULL, n.out = 53L,
  ii.match.mask = NULL, max.match.length = NULL,
  n.unpinned.observations = 0L, model = "Empirical Bayes")
```

Arguments

parent	
max.n.sims	the number of simulated curves to generate in a forecast
peak.time.dist	the distribution of smoothed-curve peak times that the prior should follow. If enabled, each smoothed curve will be x-shifted to have a peak time which is drawn from this distribution. The default setting is to disable this transformation. The default enabled distribution is a discrete uniform distribution fitted to the peak times of the smoothed curves provided to eb.createForecasts in the argument <code>fit.obj</code> .
x.shift.dist	the distribution of x-shifts to apply (after any x-shift from <code>peak.time.dist</code>). The default setting is to enable this transformation. The default enabled distribution is a discrete uniform distribution centered at zero with width equal to twice the bin width of a histogram of the <code>fit.obj</code> peak times, using Sturges' rule.
x.scale.dist	
y.scale.baseline	a single numeric value. Any y-scale transforms will only transform about and above this baseline value; for example, for a baseline of 4 and scaling factor of 2, the y-value 1 will not be scaled, since $1 < 4$, and the y-value 5 will be scaled to $5 + (5 - 4) * 2 = 6$. The default is 0, which, for non-negative smoothed curves, corresponds to simply multiplying y-values by the scaling factor.

peak.height.dist	the distribution of smoothed-curve peak heights that the prior should follow. If enabled, each smoothed curve will be y-scaled to have a peak height which is drawn from this distribution. The default setting is to disable this transformation. The default enabled distribution is a uniform distribution fitted to the <code>fit.obj</code> peak heights. If a smoothed curve remains completely below <code>y.scale.baseline</code> the entire time, then y-scaling will have no effect on that curve, and the peak height will remain at its original value. If a peak height selected from the distribution is lower than <code>y.scale.baseline</code> , parts of the curve above the baseline will be scaled by a negative factor so that the original peak is mapped to the drawn peak height value; however, this inversion is likely undesirable, and the resulting peak height may be any value between the drawn peak height value and the baseline value.
y.scale.dist	the distribution of y-scales to apply (after any y-scale from <code>peak.height.dist</code>). The default setting is to enable this transformation. The default enabled transformation is a log-uniform distribution centered at 0 in the log-scale with log-scale width equal to twice the bin width of a histogram of the logarithms of the <code>fit.obj</code> peak heights, using Sturges' rule. Note that this default behavior can significantly bias the mean of the prior for the peak heights, but does not significantly affect the median of the prior for the peak heights.
sd.option	one of "match", "scale", or "prior", which controls the assignment of a single noise level, or distribution of possible noise levels, to a transformed curve: "match" chooses the <code>sigma.hat</code> associated with the selected smooth curve; "scale" does the same, but scales this <code>sigma.hat</code> by the y-scale factor given by the transformations selected from <code>peak.height.dist</code> and <code>y.scale.dist</code> ; "prior" selects a noise level uniformly from the <code>sigma.hat</code> 's of all smoothed curves fed into the EB method, not just the one corresponding to the current transformed curve.
sd.prior	controls the distribution of noise levels used when <code>sd.option</code> is "prior"; currently, the only choice is "uniform".
sd.scale.dist	controls the distribution of noise level scaling factors, which are applied after <code>sd.option</code> and <code>sd.prior</code> are used to select an initial noise level. The default setting is to disable this transformation. There is no default enabled distribution.
reasonable.future.weight	controls the coefficient of the "reasonable future" term added to the conditional log-likelihood of the observed values given a transformed curve and noise level when calculating importance weights. The default value is 0, which disables this feature.
n.future.neighbors	controls the number of neighbors to use when determining the "reasonable future" term: for a given transformed curve and noise parameter, the neighbors are the <code>n.future.neighbors</code> historical noisy curves from the <code>dat</code> argument in eb.createForecasts with the highest log-likelihoods in future weeks (after <code>time.of.forecast</code> ; the reasonable future term is the average across these neighbors of the log-likelihood in future weeks.
bias_peaktime_weight	is the coefficient of a "peak time bias" term added to the conditional log-likelihood of the observed values given a transformed curve and noise level when calcu-

	lating importance weights. The default value is NA, which disables this feature. When it is enabled, the peak time bias term is the log-likelihood of the underlying peak time (before adding noise) of a curve, according to a Gaussian distribution with mean <code>bias_peaktime_mean</code> and standard deviation <code>bias_peaktime_sd</code> , multiplied by <code>bias_peaktime_weight</code> .
<code>bias_peaktime_mean</code>	the mean of the peak time bias Gaussian distribution
<code>bias_peaktime_sd</code>	the standard deviation of the peak time bias distribution
<code>bias_peakheight_weight</code>	is the coefficient of a "peak height bias" term added to the conditional log-likelihood of the observed values given a transformed curve and noise level when calculating importance weights. The default value is NA, which disables this feature. When it is enabled, the peak height bias term is the log-likelihood of the underlying peak height (before adding noise) of a curve, according to a Gaussian distribution with mean <code>bias_peakheight_mean</code> and standard deviation <code>bias_peakheight_sd</code> , multiplied by <code>bias_peakheight_weight</code> .
<code>bias_peakheight_mean</code>	the mean of the peak height bias Gaussian distribution
<code>bias_peakheight_sd</code>	the standard deviation of the peak height bias distribution
<code>inactive.seasons</code>	is currently ignored.
<code>n.out</code>	is the number of observations that each outputted noisy curve should contain; it should be less than or equal to the length of the shortest smooth curve. For weekly data and year-long seasons, this should be 52 or 53.
<code>ii.match.mask</code>	is a vector of indices in <code>seq_len(n.out)</code> ; only observations at these times will be considered when computing the likelihood of observations and assigning "reasonable future" terms to transformed curves and noise levels.
<code>max.match.length</code>	is a single integer controlling the maximum number of observations to use when computing the log-likelihood of <code>new.dat</code> given a transformed curve and noise level; if more than <code>max.match.length</code> observations are available at <code>time.of.forecast</code> after applying <code>ii.match.mask</code> , only the <code>max.match.length</code> most recent observations are used in the likelihood calculation.
<code>n.unpinned.observations</code>	is a single integer controlling what values in the noisy transformed curves in the posterior are "pinned" to the observed values in <code>new.dat</code> ; any observations after <code>time.of.forecast-n.unpinned.observations</code> are not pinned.

Details

Most settings are single integers, single reals, or character vectors where the first entry holds the desired value. Transformation distribution settings, on the other hand, can be one of several options, with the following associated meanings:

- NULL: Use the default setting for this transformation, either disabling the transformation or using the default distribution

- TRUE: Enable this transformation, and use the default distribution for this transformation
- FALSE: Disable this transformation
- Single integer: Enable this transformation and use the default distribution for this transformation, but break the distribution into the specified number of bins (rather than the default) when applying the grid importance sampling algorithm used by [eb.createForecasts](#)
- Distribution with bins:
- Function from curve.obj to distribution with bins:

Value

a list of parameter settings used by [eb.createForecasts](#)

Author(s)

Logan C. Brooks, David C. Farrow, Sangwon Hyun, Ryan J. Tibshirani, Roni Rosenfeld

Examples

```
default.control.list get_eb_control_list()
with.less.sims = get_eb_control_list(max.n.sims = 10000L)
with.less.sims.another.way = get_eb_control_list(default.control.list, max.n.sims = 10000L)
with.less.sims.and.sd.option.scale = get_eb_control_list(with.less.sims, sd.option="scale")
```

match.dist

Returns a distribution that can be divided into buckets — a named list containing ‘n’, a single integer representing the number of buckets into which the distribution can be broken down; ‘choices’, an n-length vector containing a representative element from each bucket; ‘probs’, an n-length vector containing the probability mass assigned to each bucket, and ‘sampler’, a function from vectors of bucket indices (in seq_len(n)) to randomly sampled elements within the corresponding buckets. Designed for input processing and validation within another function.

Description

Returns a distribution that can be divided into buckets — a named list containing ‘n’, a single integer representing the number of buckets into which the distribution can be broken down; ‘choices’, an n-length vector containing a representative element from each bucket; ‘probs’, an n-length vector containing the probability mass assigned to each bucket, and ‘sampler’, a function from vectors of bucket indices (in seq_len(n)) to randomly sampled elements within the corresponding buckets. Designed for input processing and validation within another function.

Usage

```
match.dist(curve.obj, dist, null.replacement, true.replacement,
  false.replacement, integer.replacement.fn)
```

Arguments

<code>curve.obj</code>	output of eb.fitSmoothCurves
<code>dist</code>	one of the following: (a) NULL, (b) a single boolean, (c) a single integer, (d) a function that outputs a distribution given a <code>curve.obj</code> , or (e) a distribution.
<code>null.replacement</code>	one of (b)–(e), which is used to replace NULL inputs
<code>true.replacement</code>	one of (c)–(e), which is used to replace TRUE inputs
<code>false.replacement</code>	one of (c)–(e), which is used to replace FALSE inputs
<code>integer.replacement.fn</code>	a function from a single integer to either (d) or (e), called on integer inputs to produce a replacement value

Value

a distribution, which incorporates any fitting procedure to `curve.obj`;

<code>match.single.integer</code>	<i>Returns a possibly-named length-1 non-NA integer-class vector version of the input, or produces an error if the input seems inappropriate. Designed for input processing and validation within another function.</i>
-----------------------------------	---

Description

Returns a possibly-named length-1 non-NA integer-class vector version of the input, or produces an error if the input seems inappropriate. Designed for input processing and validation within another function.

Usage

```
match.single.integer(n)
```

<code>unifChoicePrior</code>	<i>Creates a uniform distribution over discrete choices which can be used with get.eb.control.list.</i>
------------------------------	---

Description

Creates a uniform distribution over discrete choices which can be used with [get.eb.control.list](#).

Usage

```
unifChoicePrior(choices)
```

Arguments

choices a vector of (discrete) choices

Value

a uniform discrete distribution over choices.

Examples

```
uniform.seq = unifChoicePrior(letters[1:5])

## The distributions used by EB can be broken down into buckets;
## for the uniform discrete distribution, each bucket corresponds
## (boringly) to a single choice from =choices=. However, it is
## important to have a common interface.
random.bucket.indices = sample(seq_len(uniform.seq$n), 10000, replace=TRUE, prob=uniform.seq$probs)
random.elements = uniform.seq$sampler(random.bucket.indices)
random.elements.another.way = uniform.seq$choices[random.bucket.indices] # only works for =unifChoicePrior=
random.elements.a.third.way = letters[random.bucket.indices] # only works for this example
```

Index

eb.createForecasts, [2](#), [3](#), [4](#), [6](#), [7](#), [9](#)
eb.fitSmoothCurves, [3](#), [10](#)
eb.sim, [4](#)

fit.eb.control.list, [5](#)

get.eb.control.list, [10](#)
get_eb_control_list, [6](#)

match.dist, [9](#)
match.single.integer, [10](#)

unifChoicePrior, [10](#)