

Barcelona Supercomputing Center Centro Nacional de Supercomputación

BSC

# **R** user meeting

Victòria A., Ariadna B., Theertha K.

## Agenda

- 1. Ice-breaker: package 'beepr'
- 2. News
  - $\circ$  General
  - startR
  - multiApply
  - CSTools
  - CSDownscale
  - $\circ$  esviz
  - SUNSET
- 3. User presentation: BSC-ES Infrastructure (Victòria)
- 4. Q&A

# Ice-breaker: Package 'beepr'



## **Easily Play Notification Sounds on any Platform**

Package 'beepr' : The main function of this package is beep(). It is intended to play notification sounds on whatever platform you are on, for example, if you are running a long analysis in the background and want to know when it is ready.

There are a few number of different sounds to choose from.

| 1. "ping"      |                  |   |
|----------------|------------------|---|
| 2. "coin"      | L7 # ·           | install.packages("beepr")   |
| 3. "fanfare"   | 18               |   |
| 4. "complete"  | 19 #<br>20 lil   | load package<br>prary("beepr")  |
| 5. "treasure"  | 21               | and the second |
| 6. "ready"     | 22 #  <br>23 bee | olay custom sound<br>ep(8)  |
| 7. "shotgun"   | 24               |   |
| 8. "mario"     | 25 #  <br>26 bee | olay a random sound<br>ep(0)  |
| 9. "wilhelm"   | 27               |   |
| 10. "facebook" | 28 #  <br>29 bee | play custom sound on error<br>ep_on_error(1 + "a", sound = "wilhelm")   |
| 11. "sword"    | 30               |   |

# General



### **RStudio and Jupyter Notebook in Hub**

Several IDE options now available in the new BSC-ES Hubs:

- RStudio in the Hub: <u>https://earth.bsc.es/wiki/doku.php?id=computing:bsceshub#rstudio-server\_on\_t</u> <u>he\_hub</u>
- ★ Jupyter Notebooks with R kernel in the Hub: <u>https://earth.bsc.es/wiki/doku.php?id=computing:bsceshub#using\_jupyter\_note\_books</u>

Reminder of the new R modules in the hub:

R/4.4.1-gfbf-2023b (+ R-bundle-CRAN/2024.06-foss-2023b for most CRAN Packages; and R-bundle-Bioconductor/3.19-foss-2023b-R-4.4.1 for Bioconductor packages)

# startR



New startR version 2.5.0 is installed in the WS, Hub and HPC machines. CRAN submission is pending due to an issue with a dependency.

The main features of this new release are:

- ★ New special setup to use Compute() on Nord4.
- ★ Improvements to use Compute() with Autosubmit more smoothly, making the experiment creation step more flexible and consistent.
- ★ Support for conda environments when using Compute().
- ★ Bugfix in Start() to correctly round latitude and longitude values in netCDF files.

You can find the full list of changes and all details on GitLab.

# multiApply



### Apply(): Matching names in input 'data' to 'fun'

Apply() applies a function 'fun' to a list of arrays 'data'. In the current version of the multiApply(), the elements in 'data' are expected to have the same order as the arguments in 'fun'. List names are not matched to the arguments:

```
fun <- function(x, y, z) {</pre>
  return(x + y - z)
}
# Correct
res <- Apply(data = list(x, y, z),</pre>
             fun = fun)$output1
# Also correct
res <- Apply(data = list(x = x, y = y, z = z),
             fun = fun)$output1
# Technically also correct, list names are ignored!
res <- Apply(data = list(x = x, z = y, y = z),
             fun = fun)$output1
# Wrong order, will return wrong result
res <- Apply(data = list(x = x, z = z, y = y),
             fun = fun)$output1
```

### Apply(): Matching names in input 'data' to 'fun'

The new fix in Apply() matches the list names in 'data' to the arguments of 'fun' and reorders data accordingly:

issue: https://earth.bsc.es/gitlab/ces/multiApply/-/issues/20
status: in branch dev-reorder\_named\_data

### Apply(): Matching names in input 'data' to 'fun'

Functions with the `...` parameter can also be used, by keeping the corresponding elements unnamed:

issue: <a href="https://earth.bsc.es/gitlab/ces/multiApply/-/issues/20">https://earth.bsc.es/gitlab/ces/multiApply/-/issues/20</a><br/>status: in branch dev-reorder\_named\_data

# easyNCDF



The easyNCDF package provides wrappers for the ncdf4 package to simplify reading from and writing to NetCDF files, offering tools to handle multi-dimensional R arrays.

Minor release updates:

- Documentation fix in ArrayToNc() function (takes multi-dimensional R arrays and stores them in a NetCDF file).
- Maintainer change from An-Chi Ho to Ariadna Batalla.



# s2dv



### **Difference between ProbBins() and GetProbs()**

ProbBins() and GetProbs() are two s2dv functions for computing categorical forecast information.

They have a similar use and it can be unclear what the specific differences are:

### **Difference between ProbBins() and GetProbs()**

1

| Feature          | ProbBins()   | GetProbs()  |
|------------------|--|---|
| Input data       | Typically anomalies  | Forecasts/observations  |
| Ouptut           | One-hot assignments to the defined categories. Member dimension. | Probabilities for each defined category (e.g., terciles). No member dimension |
| Threshold type   | Quantile or absolute   | Quantile or absolute  |
| Cross-validation | Yes  | Yes   |
| Weights          | No   | Yes   |
| Flexibility      | Moderate   | High (reference period, bin dim, etc.)  |

# **CSTools**



## New functions BindDim() and CST\_BindDim()

New function BindDim() is a **wrapper of abind::abind()**, designed to work with arrays with named dimensions in a more user-friendly way:

```
# 10 10 3
```



## New functions BindDim() and CST\_BindDim()

CST\_BindDim() serves the same function as BindDim() for s2dv\_cubes, binding data and metadata to keep the result consistent. It has the following parameters:

- **x**: Two or more objects of class s2dv\_cube to be bound together.
- **along**: Name of the binding dimension.
- dat\_dim: Name of the dataset dimension in the s2dv\_cube. Default value is NULL.
   Specifying this dimension ensures the dataset metadata is correctly preserved.
- var\_dim: Name of the dataset dimension in the s2dv\_cube. Default value is NULL.
   Specifying this dimension ensures the dataset metadata is correctly preserved.

### MR: <u>https://earth.bsc.es/gitlab/external/cstools/-/merge\_requests/215</u> status: in master



## Implemented EvalTrainIndices in QuantileMapping

- EvalTrainIndices() generates Training and Evaluation Indices for different cross-validation methods.
- Implemented EvalTrainIndices() in QuantileMapping() to be able to use more cross-validation methods (through parameter eval.method)
  - 7 #' @param eval.method Character. The cross-validation method. Options include:
  - 8 #' -\code{leave-k-out}: Leaves out \code{k} points at a time for evaluation.
  - 9 #' -\code{retrospective}: Uses past data for training and a future point for evaluation.
  - 10 #' -\code{in-sample}: Uses the entire dataset for both training and evaluation.
  - 11 #' -\code{hindcast-vs-forecast}: Uses all years from the hindcast sdate dimension
  - 12 #' as training and all years from the forecast sdate dimension will be corrected.

# issue: <a href="https://earth.bsc.es/gitlab/external/cstools/-/issues/161">https://earth.bsc.es/gitlab/external/cstools/-/issues/161</a> status: in branch dev-quantilemapping



# **CSDownscale**



### Added examples for functions, in branch - example\_update

• Examples generally shows :

- simple parameter usage
- downscaling forecast
- downscaling using Large-scale variables (Analogs())
- downscaling for daily/monthly data
- point downscaling

path : <a href="https://earth.bsc.es/gitlab/es/csdownscale/-/tree/example\_update/inst/examples/example\_theertha">https://earth.bsc.es/gitlab/es/csdownscale/-/tree/example\_update/inst/examples/example\_theertha</a>
<br/>
status: in branch example\_updates



# esviz



### esviz:

- Contains different **plotting functions** (for maps, time series, scorecards, etc.
- Currently being prepared for its first package release, and is not yet available on CRAN.

Because these functions were originally created in different context, similar parameters sometimes had different names



We've been **standardizing** these inconsistencies to make the future package more user-friendly.



These updates are **backward compatible**: older code will still work, but warning messages will appear to notify you of the new names.

### **Parameter name standarization**

Some examples:

- `var`  $\rightarrow$  `data`
- `sizetit` and `title\_size` → `title\_scale`
- $dots_size \rightarrow dot_size$
- `dots\_shape` and `dots\_symbol` → `dot\_symbol`
- `leg` and `legend`  $\rightarrow$  `drawleg`
- ...

Please **open an issue** if you find any problems related to the name changes or if you have any suggestions/observations!

path: https://earth.bsc.es/gitlab/es/esviz/-/merge\_requests/31 status: in branch dev-unify\_esviz\_parameters

# **SUNSET**



### New full cross-validation workflow: Anomalies

CST is developing a series of workflows to allow climate forecast products to be computed using cross-validation in all steps of the workflow. The first one includes **anomaly computation** and **skill metrics**, with two functions:

**Crossval\_anomalies.R:** The hcst, fcst and obs anomalies are computed in cross-validation. Limits and forecast probabilities are also computed for the categories requested in the recipe.

**Crossval\_metrics.R:** The results of Crossval\_anomalies() are used to compute the metrics requested in the Skill section of the recipe.

Example <u>recipe</u> and <u>script</u> on GitLab

MR: <u>https://earth.bsc.es/gitlab/es/sunset/-/merge\_requests/204</u> status: in master The plots showing the probability of extremes can now be masked with the corresponding RPSS as calculated by Crossval\_metrics(). If the different sets of probabilities are named and RPSS is requested, an RPSS will be returned for each of the sets with the names given in the recipe:



The Visualization module can automatically use these names to apply the correct skill score mask to each map.

MR: <u>https://earth.bsc.es/gitlab/es/sunset/-/merge\_requests/204</u> status: in master

### Visualization: Mask extreme probabilities with RPSS

ECMWF SEAS5 (v5.1) / 2 Metre Temperature Above 90% / June 2025 / Start date: 01–05–2025



NZt. 40N R 38N 10W 2W 2F 8W 6W 4W 4F -0.2 0.2 0.6 0.8 -0.8-0.6-0.40.4 -1 0

ECMWF SEAS5 (v5.1) / 2 Metre Temperature

RPSS (P90) / June / 1993-2016

MR: <u>https://earth.bsc.es/gitlab/es/sunset/-/merge\_requests/204</u> status: in master Previously, reference dataset dates were loaded based on the metadata from the model netCDF files. If the model metadata was not correct, this could cause issues, such as a mismatch between model and reference dates.

This has been improved for **seasonal monthly** and **daily** data as well as **subseasonal daily**. SUNSET now computes the time stamps needed based on the information in the recipe.

We have also added new unit tests for the generation of the expected dates.

issue: <u>https://earth.bsc.es/gitlab/es/sunset/-/issues/173</u> status: in master

Authors: Núria Pérez-Zanón and An-Chi Ho (December 2021) Updated by Victòria Agudetse (May 2025)



Aside from the data and software in our personal laptops, we all have access to common BSC infrastructure.

### We access the BSC infrastructure:

- ★ When we connect to the <u>BSC-ES Hub</u>
- $\star$  When we use the <u>workstations</u> in the office
- ★ When we connect remotely via ssh to a workstation (bscearthXXX.int.bsc.es)
  - To ssh from windows: <u>https://earth.bsc.es/wiki/doku.php?id=computing:sshwindows</u>
  - To set up passwordless ssh connection:
     <u>https://earth.bsc.es/wiki/doku.php?id=computing:sshkeyautologon</u>
- ★ When we connect to one of the servers or HPC machines in BSC (MN5, Nord4, etc.)

When we connect to the BSC infrastructure, we find several **partitions**. A disk partition, or simply 'partition', is a segment of a hard drive that is separate and independent from other segments. Each partition serves a different purpose and is accessible from different machines.



It is also possible to connect to BSC infrastructure through **servers** (physical machines), which have different uses:

### **★** bscearth000.int.bsc.es and bscearth001.int.bsc.es

- Download data
- run the automatic package tests (GitLab CI/CD, see e.g.: <u>Pipelines · Earth Sciences / s2dv · GitLab</u>)
- transfer1.bsc.es
  - Internal transfer of data, e.g. from esarchive to GPFS and vice versa.
- ★ bscesshiny01.bsc.es
  - Shiny server, hosts shiny apps.
- ★ bscesftp.bsc.es
  - Share files externally, see:

https://earth.bsc.es/wiki/doku.php?id=computing:public\_ftp

- **b**scesautosubmit01.bsc.es and bscesautosubmit02.bsc.es
  - Dedicated machines to launch workflows with the Autosubmit workflow manager.
     To be decommissioned (date TBD)

https://earth.bsc.es/wiki/doku.php?id=tools:autosubmit

A software stack is the collection of programs and modules (including the operating system, architectural layers, protocols, runtime environments, ...) that are installed in a machine.

- The software stack at BSC can be different among different machines and departments ×  $\star$ 
  - We have access to:
    - BSC software stack (not managed by CES) 0
    - BSC-ES software stack (managed by CES) 0
      - Workstations, Nord4 and CTE-AMD already using it
      - Hubs have a different software stack, more updated
      - In some machines, we should edit the **bashrc** to use it (instructions are always in the wiki: <a href="https://earth.bsc.es/wiki/doku.php?id=library:computing">https://earth.bsc.es/wiki/doku.php?id=library:computing</a>)
      - It is built on modules, some useful commands are:
        - module list # show all loaded modules
        - module load \* #load the '\*' module
        - module av \* # show all available modules matching '\*'
      - other software programs like mendeley can be opened in the WS/Hub: /shared/earth/software/mendeley/latest/bin/mendeleydesktop

Open an issue in the Requests GitLab to ask for new software or R packages (@smirosla)

What information do we need to know for each machine?

- Does it have BSC-ES software?
- is /esarchive/ mounted?
- Internet access?
- How to send jobs?
- Memory per node, cores per node....

Hub Workstations (WS) Marenostrum 5 AMD cluster Nord4

Find the information here: <u>https://earth.bsc.es/wiki/doku.php?id=library:computing</u>

On the HPC machines (Nord4, CTE-AMD, MN5) there are two types of nodes:

- ★ Login nodes: The node we are in when we to connect to the machine via ssh. It can be used to execute basic tasks and send jobs. Should not be used for computation.
- ★ Compute nodes: For computation and memory-intensive tasks. They are assigned by the scheduler (e.g. slurm) when sending jobs or requesting interactive sessions. The job scheduler assigns resources on compute nodes based on priority.

More information here: <u>https://www.bsc.es/supportkc/docs/Nord4/slurm</u> <u>https://www.bsc.es/supportkc/docs/MareNostrum5/slurm</u> <u>https://www.bsc.es/supportkc/docs/CTE-AMD/slurm</u>

#### Workstations

- R/4.1.2
- For debugging code (small data) or running startR workflows in remote machines
- Internet connection
- BSC-ES software stack
- /esarchive is mounted
- To be decommissioned

#### Hub (bsceshub03-07)

- R/4.4.1
- For debugging code (small data), startR workflows or running autosubmit.
- Internet connection
- BSC-ES software stack
- /esarchive is mounted
- Replaces workstations.

#### Nord4

- R/4.1.2
- For memory-intensive jobs
- Internet access in login node 0
- BSC-ES software stack; must use nord3\_singu\_es wrapper
- /esarchive is mounted

#### Marenostrum 5

- R/4.3.3 or R/4.2.2 (SUNSET conda)
- For memory-intensive jobs
- Some BSC-ES software stack available, conda environments can be installed
- Internet access in login node 4
- no access to /esarchive



#### CTE-AMD

- R/4.1.2 or R/4.3.3 (for R-INLA)
- For memory-intensive jobs
- BSC-ES software stack
- no access to /esarchive

#### Recommendations

#### ★ Save your scripts in GitLab (intermediate and final versions)

- In an existing GitLab project
- In a personal project
- Documentation: <u>https://earth.bsc.es/wiki/doku.php?id=library:computing#git</u>
- If you have internet connection, you can source your code directly from GitLab
- Clone repositories under /esarchive/scratch/<username>/
  - You will have internet connection to push your changes
  - The code will be accessible from Hub, Nord4 and shiny server
  - There is no back-up copy of /esarchive (another good reason to use GitLab)

#### Don't install local versions of R packages

- If you do, we cannot debug the code and reproduce the errors
- Better to open an issue in Requests to ask for the installation: it's easier to debug and everyone can use it

### Infrastructure in the wiki:

https://earth.bsc.es/wiki/doku.php?id=library:best\_practices#network\_infrastructure

Q&A: What else do we need to know? What questions do we have?



# Thanks for joining



Next meeting: TBA