

# Package ‘scoringutils’

November 29, 2023

**Title** Utilities for Scoring and Assessing Predictions

**Version** 1.2.2

**Language** en-GB

**Description** Provides a collection of metrics and proper scoring rules (Tilmann Gneiting & Adrian E Raftery (2007) <[doi:10.1198/016214506000001437](https://doi.org/10.1198/016214506000001437)>, Jordan, A., Krüger, F., & Lerch, S. (2019) <[doi:10.18637/jss.v090.i12](https://doi.org/10.18637/jss.v090.i12)>) within a consistent framework for evaluation, comparison and visualisation of forecasts. In addition to proper scoring rules, functions are provided to assess bias, sharpness and calibration (Sebastian Funk, Anton Camacho, Adam J. Kucharski, Rachel Lowe, Rosalind M. Eggo, W. John Edmunds (2019) <[doi:10.1371/journal.pcbi.1006785](https://doi.org/10.1371/journal.pcbi.1006785)>) of forecasts. Several types of predictions (e.g. binary, discrete, continuous) which may come in different formats (e.g. forecasts represented by predictive samples or by quantiles of the predictive distribution) can be evaluated. Scoring metrics can be used either through a convenient data.frame format, or can be applied as individual functions in a vector / matrix format. All functionality has been implemented with a focus on performance and is robustly tested. Find more information about the package in the accompanying paper (<[doi:10.48550/arXiv.2205.07090](https://doi.org/10.48550/arXiv.2205.07090)>).

**License** MIT + file LICENSE

**Encoding** UTF-8

**LazyData** true

**Imports** data.table, ggdist (>= 3.2.0), ggplot2 (>= 3.4.0), lifecycle, methods, rlang, scoringRules, stats

**Suggests** kableExtra, knitr, magrittr, rmarkdown, testthat, vdiff

**Config/Needs/website** r-lib/pkgdown, amirmasoudabdol/preferably

**Config/testthat/edition** 3

**RoxygenNote** 7.2.3

**URL** <https://doi.org/10.48550/arXiv.2205.07090>,  
<https://epiforecasts.io/scoringutils/>,  
<https://github.com/epiforecasts/scoringutils>

**BugReports** <https://github.com/epiforecasts/scoringutils/issues>

**VignetteBuilder** knitr

**Depends** R (>= 3.6)

**NeedsCompilation** no

**Author** Nikos Bosse [aut, cre] (<<https://orcid.org/0000-0002-7750-5280>>),  
 Sam Abbott [aut] (<<https://orcid.org/0000-0001-8057-8037>>),  
 Hugo Gruson [aut] (<<https://orcid.org/0000-0002-4094-1476>>),  
 Johannes Bracher [ctb] (<<https://orcid.org/0000-0002-3777-1410>>),  
 Sebastian Funk [aut]

**Maintainer** Nikos Bosse <[nikosbosse@gmail.com](mailto:nikosbosse@gmail.com)>

**Repository** CRAN

**Date/Publication** 2023-11-29 15:50:10 UTC

## R topics documented:

abs_error . . . . .	3
add_coverage . . . . .	4
ae_median_quantile . . . . .	5
ae_median_sample . . . . .	6
available_metrics . . . . .	6
avail_forecasts . . . . .	7
bias_quantile . . . . .	8
bias_range . . . . .	9
bias_sample . . . . .	11
brier_score . . . . .	12
check_forecasts . . . . .	13
correlation . . . . .	15
crps_sample . . . . .	15
dss_sample . . . . .	16
example_binary . . . . .	17
example_continuous . . . . .	18
example_integer . . . . .	19
example_point . . . . .	19
example_quantile . . . . .	20
example_quantile_forecasts_only . . . . .	21
example_truth_only . . . . .	22
find_duplicates . . . . .	23
interval_score . . . . .	23
logs_binary . . . . .	25
logs_sample . . . . .	26
log_shift . . . . .	27
mad_sample . . . . .	28
make_NA . . . . .	29
merge_pred_and_obs . . . . .	30
metrics . . . . .	31
pairwise_comparison . . . . .	31

pit . . . . .	33
pit_sample . . . . .	34
plot_avail_forecasts . . . . .	35
plot_correlation . . . . .	36
plot_heatmap . . . . .	37
plot_interval_coverage . . . . .	38
plot_pairwise_comparison . . . . .	38
plot_pit . . . . .	39
plot_predictions . . . . .	40
plot_quantile_coverage . . . . .	41
plot_ranges . . . . .	42
plot_score_table . . . . .	43
plot_wis . . . . .	44
print.scoringutils_check . . . . .	45
quantile_score . . . . .	46
sample_to_quantile . . . . .	47
score . . . . .	48
set_forecast_unit . . . . .	50
se_mean_sample . . . . .	51
squared_error . . . . .	52
summarise_scores . . . . .	52
theme_scoringutils . . . . .	55
transform_forecasts . . . . .	55

## Index 58

---

abs_error	<i>Absolute Error</i>
-----------	-----------------------

---

### Description

Calculate absolute error as

$$\text{abs}(\text{true\_value} - \text{median\_prediction})$$

### Usage

```
abs_error(true_values, predictions)
```

### Arguments

true_values	A vector with the true observed values of size n
predictions	numeric vector with predictions, corresponding to the quantiles in a second vector, quantiles.

### Value

vector with the absolute error

**See Also**

[ae\\_median\\_sample\(\)](#), [ae\\_median\\_quantile\(\)](#)

**Examples**

```
true_values <- rnorm(30, mean = 1:30)
predicted_values <- rnorm(30, mean = 1:30)
abs_error(true_values, predicted_values)
```

---

add_coverage	<i>Add coverage of central prediction intervals</i>
--------------	---

---

**Description**

Adds a column with the coverage of central prediction intervals to unsummarised scores as produced by [score\(\)](#)

**Usage**

```
add_coverage(scores, by, ranges = c(50, 90))
```

**Arguments**

scores	A data.table of scores as produced by <a href="#">score()</a> .
by	character vector with column names to add the coverage for.
ranges	numeric vector of the ranges of the central prediction intervals for which coverage values shall be added.

**Details**

The coverage values that are added are computed according to the values specified in by. If, for example, by = "model", then there will be one coverage value for every model and [add\\_coverage\(\)](#) will compute the coverage for every model across the values present in all other columns which define the unit of a single forecast.

**Value**

a data.table with unsummarised scores with columns added for the coverage of the central prediction intervals. While the overall data.table is still unsummarised, note that for the coverage columns some level of summary is present according to the value specified in by.

**Examples**

```
library(magrittr) # pipe operator

score(example_quantile) %>%
  add_coverage(by = c("model", "target_type")) %>%
  summarise_scores(by = c("model", "target_type")) %>%
  summarise_scores(fun = signif, digits = 2)
```

---

ae\_median\_quantile     *Absolute Error of the Median (Quantile-based Version)*

---

## Description

Absolute error of the median calculated as

$$\text{abs}(\text{true\_value} - \text{prediction})$$

The function was created for internal use within `score()`, but can also be used as a standalone function.

## Usage

```
ae_median_quantile(true_values, predictions, quantiles = NULL)
```

## Arguments

<code>true_values</code>	A vector with the true observed values of size <code>n</code>
<code>predictions</code>	numeric vector with predictions, corresponding to the quantiles in a second vector, <code>quantiles</code> .
<code>quantiles</code>	numeric vector that denotes the quantile for the values in predictions. Only those predictions where <code>quantiles == 0.5</code> will be kept. If <code>quantiles</code> is <code>NULL</code> , then all predictions and <code>true_values</code> will be used (this is then the same as <code>abs_error()</code> )

## Value

vector with the scoring values

## See Also

`ae_median_sample()`, `abs_error()`

## Examples

```
true_values <- rnorm(30, mean = 1:30)
predicted_values <- rnorm(30, mean = 1:30)
ae_median_quantile(true_values, predicted_values, quantiles = 0.5)
```

ae\_median\_sample      *Absolute Error of the Median (Sample-based Version)*

---

**Description**

Absolute error of the median calculated as

$$\text{abs}(\text{true\_value} - \text{median\_prediction})$$
**Usage**

```
ae_median_sample(true_values, predictions)
```

**Arguments**

`true_values`      A vector with the true observed values of size `n`  
`predictions`      `nxN` matrix of predictive samples, `n` (number of rows) being the number of data points and `N` (number of columns) the number of Monte Carlo samples. Alternatively, `predictions` can just be a vector of size `n`.

**Value**

vector with the scoring values

**See Also**

[ae\\_median\\_quantile\(\)](#), [abs\\_error\(\)](#)

**Examples**

```
true_values <- rnorm(30, mean = 1:30)
predicted_values <- rnorm(30, mean = 1:30)
ae_median_sample(true_values, predicted_values)
```

---

available\_metrics      *Available metrics in scoringutils*

---

**Description**

Available metrics in scoringutils

**Usage**

```
available_metrics()
```

**Value**

A vector with the name of all available metrics

---

avail_forecasts	<i>Display Number of Forecasts Available</i>
-----------------	--

---

## Description

Given a data set with forecasts, count the number of available forecasts for arbitrary grouping (e.g. the number of forecasts per model, or the number of forecasts per model and location). This is useful to determine whether there are any missing forecasts.

## Usage

```
avail_forecasts(data, by = NULL, collapse = c("quantile", "sample"))
```

## Arguments

data	<p>A <code>data.frame</code> or <code>data.table</code> with the predictions and observations. For scoring using <code>score()</code>, the following columns need to be present:</p> <ul style="list-style-type: none"> <li>• <code>true_value</code> - the true observed values</li> <li>• <code>prediction</code> - predictions or predictive samples for one true value. (You only don't need to provide a prediction column if you want to score quantile forecasts in a wide range format.)</li> </ul> <p>For scoring integer and continuous forecasts a <code>sample</code> column is needed:</p> <ul style="list-style-type: none"> <li>• <code>sample</code> - an index to identify the predictive samples in the prediction column generated by one model for one true value. Only necessary for continuous and integer forecasts, not for binary predictions.</li> </ul> <p>For scoring predictions in a quantile-format forecast you should provide a column called <code>quantile</code>:</p> <ul style="list-style-type: none"> <li>• <code>quantile</code>: quantile to which the prediction corresponds</li> </ul> <p>In addition a <code>model</code> column is suggested and if not present this will be flagged and added to the input data with all forecasts assigned as an "unspecified model". You can check the format of your data using <code>check_forecasts()</code> and there are examples for each format (<a href="#">example_quantile</a>, <a href="#">example_continuous</a>, <a href="#">example_integer</a>, and <a href="#">example_binary</a>).</p>
by	<p>character vector or <code>NULL</code> (the default) that denotes the categories over which the number of forecasts should be counted. By default (<code>by = NULL</code>) this will be the unit of a single forecast (i.e. all available columns (apart from a few "protected" columns such as 'prediction' and 'true value') plus "quantile" or "sample" where present).</p>
collapse	<p>character vector (default is <code>c("quantile", "sample")</code>) with names of categories for which the number of rows should be collapsed to one when counting. For example, a single forecast is usually represented by a set of several quantiles or samples and collapsing these to one makes sure that a single forecast only gets counted once.</p>

**Value**

A data.table with columns as specified in by and an additional column with the number of forecasts.

**Examples**

```
avail_forecasts(example_quantile,
  collapse = c("quantile"),
  by = c("model", "target_type")
)
```

---

bias_quantile	<i>Determines Bias of Quantile Forecasts</i>
---------------	--

---

**Description**

Determines bias from quantile forecasts. For an increasing number of quantiles this measure converges against the sample based bias version for integer and continuous forecasts.

**Usage**

```
bias_quantile(predictions, quantiles, true_value)
```

**Arguments**

predictions	vector of length corresponding to the number of quantiles that holds predictions
quantiles	vector of corresponding size with the quantiles for which predictions were made. If this does not contain the median (0.5) then the median is imputed as being the mean of the two innermost quantiles.
true_value	a single true value

**Details**

For quantile forecasts, bias is measured as

$$B_t = (1 - 2 \cdot \max\{i | q_{t,i} \in Q_t \wedge q_{t,i} \leq x_t\}) \mathbf{1}(x_t \leq q_{t,0.5}) + (1 - 2 \cdot \min\{i | q_{t,i} \in Q_t \wedge q_{t,i} \geq x_t\}) \mathbf{1}(x_t \geq q_{t,0.5}),$$

where  $Q_t$  is the set of quantiles that form the predictive distribution at time  $t$ . They represent our belief about what the true value  $x_t$  will be. For consistency, we define  $Q_t$  such that it always includes the element  $q_{t,0} = -\infty$  and  $q_{t,1} = \infty$ .  $\mathbf{1}(\cdot)$  is the indicator function that is 1 if the condition is satisfied and 0 otherwise. In clearer terms,  $B_t$  is defined as the maximum percentile rank for which the corresponding quantile is still below the true value, if the true value is smaller than the median of the predictive distribution. If the true value is above the median of the predictive distribution, then  $B_t$  is the minimum percentile rank for which the corresponding quantile is still larger than the true value. If the true value is exactly the median, both terms cancel out and



$B_t$  is zero. For a large enough number of quantiles, the percentile rank will equal the proportion of predictive samples below the observed true value, and this metric coincides with the one for continuous forecasts.

Bias can assume values between -1 and 1 and is 0 ideally (i.e. unbiased).

**Value**

scalar with the quantile bias for a single quantile prediction

**Author(s)**

Nikos Bosse <nikosbosse@gmail.com>

**Examples**

```
predictions <- c(
  705.500, 1127.000, 4006.250, 4341.500, 4709.000, 4821.996,
  5340.500, 5451.000, 5703.500, 6087.014, 6329.500, 6341.000,
  6352.500, 6594.986, 6978.500, 7231.000, 7341.500, 7860.004,
  7973.000, 8340.500, 8675.750, 11555.000, 11976.500
)

quantiles <- c(0.01, 0.025, seq(0.05, 0.95, 0.05), 0.975, 0.99)

true_value <- 8062

bias_quantile(predictions, quantiles, true_value = true_value)
```

---

bias_range	<i>Determines Bias of Quantile Forecasts based on the range of the prediction intervals</i>
------------	---

---

**Description**

Determines bias from quantile forecasts based on the range of the prediction intervals. For an increasing number of quantiles this measure converges against the sample based bias version for integer and continuous forecasts.

**Usage**

```
bias_range(lower, upper, range, true_value)
```

**Arguments**

lower	vector of length corresponding to the number of central prediction intervals that holds predictions for the lower bounds of a prediction interval
upper	vector of length corresponding to the number of central prediction intervals that holds predictions for the upper bounds of a prediction interval
range	vector of corresponding size with information about the width of the central prediction interval
true_value	a single true value

**Details**

For quantile forecasts, bias is measured as

$$B_t = (1 - 2 \cdot \max\{i \mid q_{t,i} \in Q_t \wedge q_{t,i} \leq x_t\}) \mathbf{1}(x_t \leq q_{t,0.5}) + (1 - 2 \cdot \min\{i \mid q_{t,i} \in Q_t \wedge q_{t,i} \geq x_t\}) \mathbf{1}(x_t \geq q_{t,0.5}),$$

where  $Q_t$  is the set of quantiles that form the predictive distribution at time  $t$ . They represent our belief about what the true value  $x_t$  will be. For consistency, we define  $Q_t$  such that it always includes the element  $q_{t,0} = -\infty$  and  $q_{t,1} = \infty$ .  $\mathbf{1}(\cdot)$  is the indicator function that is 1 if the condition is satisfied and 0 otherwise. In clearer terms,  $B_t$  is defined as the maximum percentile rank for which the corresponding quantile is still below the true value, if the true value is smaller than the median of the predictive distribution. If the true value is above the median of the predictive distribution, then  $B_t$  is the minimum percentile rank for which the corresponding quantile is still larger than the true value. If the true value is exactly the median, both terms cancel out and  $B_t$  is zero. For a large enough number of quantiles, the percentile rank will equal the proportion of predictive samples below the observed true value, and this metric coincides with the one for continuous forecasts.

Bias can assume values between -1 and 1 and is 0 ideally.

**Value**

scalar with the quantile bias for a single quantile prediction

**Author(s)**

Nikos Bosse <nikosbosse@gmail.com>

**See Also**

bias\_quantile bias\_sample

**Examples**

```
lower <- c(
  6341.000, 6329.500, 6087.014, 5703.500,
  5451.000, 5340.500, 4821.996, 4709.000,
  4341.500, 4006.250, 1127.000, 705.500
```

```

)
upper <- c(
  6341.000, 6352.500, 6594.986, 6978.500,
  7231.000, 7341.500, 7860.004, 7973.000,
  8340.500, 8675.750, 11555.000, 11976.500
)

range <- c(0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 95, 98)

true_value <- 8062

bias_range(
  lower = lower, upper = upper,
  range = range, true_value = true_value
)

```

---

bias_sample	<i>Determines bias of forecasts</i>
-------------	-------------------------------------

---

### Description

Determines bias from predictive Monte-Carlo samples. The function automatically recognises, whether forecasts are continuous or integer valued and adapts the Bias function accordingly.

### Usage

```
bias_sample(true_values, predictions)
```

### Arguments

true_values	A vector with the true observed values of size n
predictions	nxN matrix of predictive samples, n (number of rows) being the number of data points and N (number of columns) the number of Monte Carlo samples. Alternatively, predictions can just be a vector of size n.

### Details

For continuous forecasts, Bias is measured as

$$B_t(P_t, x_t) = 1 - 2 * (P_t(x_t))$$

where  $P_t$  is the empirical cumulative distribution function of the prediction for the true value  $x_t$ . Computationally,  $P_t(x_t)$  is just calculated as the fraction of predictive samples for  $x_t$  that are smaller than  $x_t$ .

For integer valued forecasts, Bias is measured as

$$B_t(P_t, x_t) = 1 - (P_t(x_t) + P_t(x_t + 1))$$

to adjust for the integer nature of the forecasts.

In both cases, Bias can assume values between -1 and 1 and is 0 ideally.

### Value

vector of length n with the biases of the predictive samples with respect to the true values.

### Author(s)

Nikos Bosse <nikosbosse@gmail.com>

### References

The integer valued Bias function is discussed in Assessing the performance of real-time epidemic forecasts: A case study of Ebola in the Western Area region of Sierra Leone, 2014-15 Funk S, Camacho A, Kucharski AJ, Lowe R, Eggo RM, et al. (2019) Assessing the performance of real-time epidemic forecasts: A case study of Ebola in the Western Area region of Sierra Leone, 2014-15. PLOS Computational Biology 15(2): e1006785. doi:10.1371/journal.pcbi.1006785

### Examples

```
## integer valued forecasts
true_values <- rpois(30, lambda = 1:30)
predictions <- replicate(200, rpois(n = 30, lambda = 1:30))
bias_sample(true_values, predictions)

## continuous forecasts
true_values <- rnorm(30, mean = 1:30)
predictions <- replicate(200, rnorm(30, mean = 1:30))
bias_sample(true_values, predictions)
```

---

brier\_score

*Brier Score*

---

### Description

Computes the Brier Score for probabilistic forecasts of binary outcomes.

### Usage

```
brier_score(true_values, predictions)
```

### Arguments

true\_values     A vector with the true observed values of size n with all values equal to either 0 or 1

predictions     A vector with a predicted probability that true\_value = 1.

**Details**

The Brier score is a proper score rule that assesses the accuracy of probabilistic binary predictions. The outcomes can be either 0 or 1, the predictions must be a probability that the true outcome will be 1.

The Brier Score is then computed as the mean squared error between the probabilistic prediction and the true outcome.

$$\text{Brier\_Score} = \frac{1}{N} \sum_{t=1}^n (\text{prediction}_t - \text{outcome}_t)^2$$

**Value**

A numeric value with the Brier Score, i.e. the mean squared error of the given probability forecasts

**Examples**

```
true_values <- sample(c(0, 1), size = 30, replace = TRUE)
predictions <- runif(n = 30, min = 0, max = 1)

brier_score(true_values, predictions)
```

---

check\_forecasts

*Check forecasts*

---

**Description**

Function to check the input data before running `score()`.

The data should come in one of three different formats:

- A format for binary predictions (see [example\\_binary](#))
- A sample-based format for discrete or continuous predictions (see [example\\_continuous](#) and [example\\_integer](#))
- A quantile-based format (see [example\\_quantile](#))

**Usage**

```
check_forecasts(data)
```

**Arguments**

**data** A data.frame or data.table with the predictions and observations. For scoring using `score()`, the following columns need to be present:

- `true_value` - the true observed values
- `prediction` - predictions or predictive samples for one true value. (You only don't need to provide a prediction column if you want to score quantile forecasts in a wide range format.)

For scoring integer and continuous forecasts a `sample` column is needed:

- `sample` - an index to identify the predictive samples in the prediction column generated by one model for one true value. Only necessary for continuous and integer forecasts, not for binary predictions.

For scoring predictions in a quantile-format forecast you should provide a column called `quantile`:

- `quantile`: quantile to which the prediction corresponds

In addition a `model` column is suggested and if not present this will be flagged and added to the input data with all forecasts assigned as an "unspecified model").

You can check the format of your data using `check_forecasts()` and there are examples for each format ([example\\_quantile](#), [example\\_continuous](#), [example\\_integer](#), and [example\\_binary](#)).

## Value

A list with elements that give information about what `scoringutils` thinks you are trying to do and potential issues.

- `target_type` the type of the prediction target as inferred from the input: 'binary', if all values in `true_value` are either 0 or 1 and values in `prediction` are between 0 and 1, 'discrete' if all true values are integers. and 'continuous' if not.
- `prediction_type` inferred type of the prediction. 'quantile', if there is a column called 'quantile', else 'discrete' if all values in `prediction` are integer, else 'continuous'.
- `forecast_unit` unit of a single forecast, i.e. the grouping that uniquely defines a single forecast. This is assumed to be all present columns apart from the following protected columns: `c("prediction", "true_value", "sample", "quantile", "range", "boundary")`. It is important that you remove all unnecessary columns before scoring.
- `rows_per_forecast` a data.frame that shows how many rows (usually quantiles or samples there are available per forecast. If a forecast model has several entries, then there a forecasts with differing numbers of quantiles / samples.
- `unique_values` A data.frame that shows how many unique values there are present per model and column in the data. This doesn't directly show missing values, but rather the maximum number of unique values across the whole data.
- `warnings` A vector with warnings. These can be ignored if you know what you are doing.
- `errors` A vector with issues that will cause an error when running `score()`.
- `messages` A verbal explanation of the information provided above.

## Author(s)

Nikos Bosse <[nikosbosse@gmail.com](mailto:nikosbosse@gmail.com)>

## See Also

Function to move from sample-based to quantile format: [sample\\_to\\_quantile\(\)](#)

**Examples**

```
check <- check_forecasts(example_quantile)
print(check)
check_forecasts(example_binary)
```

---

correlation	<i>Correlation Between Metrics</i>
-------------	------------------------------------

---

**Description**

Calculate the correlation between different metrics for a data.frame of scores as produced by [score\(\)](#).

**Usage**

```
correlation(scores, metrics = NULL)
```

**Arguments**

scores	A data.table of scores as produced by <a href="#">score()</a> .
metrics	A character vector with the metrics to show. If set to NULL (default), all metrics present in scores will be shown

**Value**

A data.table with correlations for the different metrics

**Examples**

```
scores <- score(example_quantile)
correlation(scores)
```

---

crps_sample	<i>Ranked Probability Score</i>
-------------	---------------------------------

---

**Description**

Wrapper around the [crps\\_sample\(\)](#) function from the **scoringRules** package. Can be used for continuous as well as integer valued forecasts

**Usage**

```
crps_sample(true_values, predictions)
```

**Arguments**

`true_values` A vector with the true observed values of size `n`

`predictions` `nxN` matrix of predictive samples, `n` (number of rows) being the number of data points and `N` (number of columns) the number of Monte Carlo samples. Alternatively, predictions can just be a vector of size `n`.

**Value**

vector with the scoring values

**References**

Alexander Jordan, Fabian Krüger, Sebastian Lerch, Evaluating Probabilistic Forecasts with scoringRules, <https://www.jstatsoft.org/article/view/v090i12>

**Examples**

```
true_values <- rpois(30, lambda = 1:30)
predictions <- replicate(200, rpois(n = 30, lambda = 1:30))
crps_sample(true_values, predictions)
```

---

dss\_sample

*Dawid-Sebastiani Score*

---

**Description**

Wrapper around the `dss_sample()` function from the **scoringRules** package.

**Usage**

```
dss_sample(true_values, predictions)
```

**Arguments**

`true_values` A vector with the true observed values of size `n`

`predictions` `nxN` matrix of predictive samples, `n` (number of rows) being the number of data points and `N` (number of columns) the number of Monte Carlo samples. Alternatively, predictions can just be a vector of size `n`.

**Value**

vector with scoring values

**References**

Alexander Jordan, Fabian Krüger, Sebastian Lerch, Evaluating Probabilistic Forecasts with scoringRules, <https://www.jstatsoft.org/article/view/v090i12>



## Examples

```
true_values <- rpois(30, lambda = 1:30)
predictions <- replicate(200, rpois(n = 30, lambda = 1:30))
dss_sample(true_values, predictions)
```

---

example\_binary

*Binary Forecast Example Data*

---

## Description

A data set with binary predictions for COVID-19 cases and deaths constructed from data submitted to the European Forecast Hub.

## Usage

```
example_binary
```

## Format

A data frame with 346 rows and 10 columns:

**location** the country for which a prediction was made

**location\_name** name of the country for which a prediction was made

**target\_end\_date** the date for which a prediction was made

**target\_type** the target to be predicted (cases or deaths)

**true\_value** true observed values

**forecast\_date** the date on which a prediction was made

**model** name of the model that generated the forecasts

**horizon** forecast horizon in weeks

**prediction** predicted value

## Details

Predictions in the data set were constructed based on the continuous example data by looking at the number of samples below the mean prediction. The outcome was constructed as whether or not the actually observed value was below or above that mean prediction. This should not be understood as sound statistical practice, but rather as a practical way to create an example data set.

The data was created using the script create-example-data.R in the inst/ folder (or the top level folder in a compiled package).

## Source

```
https://github.com/european-modelling-hubs/covid19-forecast-hub-europe/commit/a42867b1ea152c57e25b0
# nolint
```

---

example\_continuous      *Continuous Forecast Example Data*

---

### Description

A data set with continuous predictions for COVID-19 cases and deaths constructed from data submitted to the European Forecast Hub.

### Usage

example\_continuous

### Format

A data frame with 13,429 rows and 10 columns:

**location** the country for which a prediction was made

**target\_end\_date** the date for which a prediction was made

**target\_type** the target to be predicted (cases or deaths)

**true\_value** true observed values

**location\_name** name of the country for which a prediction was made

**forecast\_date** the date on which a prediction was made

**model** name of the model that generated the forecasts

**horizon** forecast horizon in weeks

**prediction** predicted value

**sample** id for the corresponding sample

### Details

The data was created using the script create-example-data.R in the inst/ folder (or the top level folder in a compiled package).

### Source

<https://github.com/european-modelling-hubs/covid19-forecast-hub-europe/commit/a42867b1ea152c57e25b0>  
# nolint

---

example_integer	<i>Integer Forecast Example Data</i>
-----------------	--------------------------------------

---

**Description**

A data set with integer predictions for COVID-19 cases and deaths constructed from data submitted to the European Forecast Hub.

**Usage**

example\_integer

**Format**

A data frame with 13,429 rows and 10 columns:

**location** the country for which a prediction was made  
**target\_end\_date** the date for which a prediction was made  
**target\_type** the target to be predicted (cases or deaths)  
**true\_value** true observed values  
**location\_name** name of the country for which a prediction was made  
**forecast\_date** the date on which a prediction was made  
**model** name of the model that generated the forecasts  
**horizon** forecast horizon in weeks  
**prediction** predicted value  
**sample** id for the corresponding sample

**Details**

The data was created using the script create-example-data.R in the inst/ folder (or the top level folder in a compiled package).

---

example_point	<i>Point Forecast Example Data</i>
---------------	------------------------------------

---

**Description**

A data set with predictions for COVID-19 cases and deaths submitted to the European Forecast Hub. This data set is like the quantile example data, only that the median has been replaced by a point forecast.

**Usage**

example\_point

**Format**

A data frame with

**location** the country for which a prediction was made  
**target\_end\_date** the date for which a prediction was made  
**target\_type** the target to be predicted (cases or deaths)  
**true\_value** true observed values  
**location\_name** name of the country for which a prediction was made  
**forecast\_date** the date on which a prediction was made  
**quantile** quantile of the corresponding prediction  
**prediction** predicted value  
**model** name of the model that generated the forecasts  
**horizon** forecast horizon in weeks

**Details**

The data was created using the script create-example-data.R in the inst/ folder (or the top level folder in a compiled package).

**Source**

<https://github.com/european-modelling-hubs/covid19-forecast-hub-europe/commit/a42867b1ea152c57e25b0>  
# nolint

---

example_quantile	<i>Quantile Example Data</i>
------------------	------------------------------

---

**Description**

A data set with predictions for COVID-19 cases and deaths submitted to the European Forecast Hub.

**Usage**

```
example_quantile
```

**Format**

A data frame with

**location** the country for which a prediction was made  
**target\_end\_date** the date for which a prediction was made  
**target\_type** the target to be predicted (cases or deaths)  
**true\_value** true observed values

**location\_name** name of the country for which a prediction was made

**forecast\_date** the date on which a prediction was made

**quantile** quantile of the corresponding prediction

**prediction** predicted value

**model** name of the model that generated the forecasts

**horizon** forecast horizon in weeks

## Details

The data was created using the script create-example-data.R in the inst/ folder (or the top level folder in a compiled package).

## Source

<https://github.com/european-modelling-hubs/covid19-forecast-hub-europe/commit/a42867b1ea152c57e25b0>  
# nolint

---

example\_quantile\_forecasts\_only

*Quantile Example Data - Forecasts only*

---

## Description

A data set with quantile predictions for COVID-19 cases and deaths submitted to the European Forecast Hub.

## Usage

```
example_quantile_forecasts_only
```

## Format

A data frame with 7,581 rows and 9 columns:

**location** the country for which a prediction was made

**target\_end\_date** the date for which a prediction was made

**target\_type** the target to be predicted (cases or deaths)

**forecast\_date** the date on which a prediction was made

**quantile** quantile of the corresponding prediction

**prediction** predicted value

**model** name of the model that generated the forecasts

**horizon** forecast horizon in weeks

### Details

The data was created using the script create-example-data.R in the inst/ folder (or the top level folder in a compiled package).

### Source

<https://github.com/european-modelling-hubs/covid19-forecast-hub-europe/commit/a42867b1ea152c57e25b0>  
# nolint

---

example_truth_only	<i>Truth data only</i>
--------------------	------------------------

---

### Description

A data set with truth values for COVID-19 cases and deaths submitted to the European Forecast Hub.

### Usage

```
example_truth_only
```

### Format

A data frame with 140 rows and 5 columns:

**location** the country for which a prediction was made

**target\_end\_date** the date for which a prediction was made

**target\_type** the target to be predicted (cases or deaths)

**true\_value** true observed values

**location\_name** name of the country for which a prediction was made

### Details

The data was created using the script create-example-data.R in the inst/ folder (or the top level folder in a compiled package).

### Source

<https://github.com/european-modelling-hubs/covid19-forecast-hub-europe/commit/a42867b1ea152c57e25b0>  
# nolint

---

find_duplicates	<i>Find duplicate forecasts</i>
-----------------	---------------------------------

---

**Description**

Helper function to identify duplicate forecasts, i.e. instances where there is more than one forecast for the same prediction target.

**Usage**

```
find_duplicates(data, forecast_unit)
```

**Arguments**

data	A data.frame as used for <code>score()</code>
forecast_unit	A character vector with the column names that define the unit of a single forecast. If missing the function tries to infer the unit of a single forecast.

**Value**

A data.frame with all rows for which a duplicate forecast was found

**Examples**

```
example <- rbind(example_quantile, example_quantile[1000:1010])
find_duplicates(example)
```

---

interval_score	<i>Interval Score</i>
----------------	-----------------------

---

**Description**

Proper Scoring Rule to score quantile predictions, following Gneiting and Raftery (2007). Smaller values are better.

The score is computed as

$$\text{score} = (\text{upper} - \text{lower}) + \frac{2}{\alpha} (\text{lower} - \text{true\_value}) * \mathbf{1}(\text{true\_value} < \text{lower}) + \frac{2}{\alpha} (\text{true\_value} - \text{upper}) * \mathbf{1}(\text{true\_value} > \text{upper})$$

where  $\mathbf{1}()$  is the indicator function and indicates how much is outside the prediction interval.  $\alpha$  is the decimal value that indicates how much is outside the prediction interval.

To improve usability, the user is asked to provide an interval range in percentage terms, i.e. `interval_range = 90` (percent) for a 90 percent prediction interval. Correspondingly, the user would have to provide the 5% and 95% quantiles (the corresponding alpha would then be 0.1). No specific distribution is assumed, but the range has to be symmetric (i.e you can't use the 0.1 quantile as the lower bound and the 0.7 quantile as the upper). Non-symmetric quantiles can be scored using the function `quantile_score()`.

**Usage**

```
interval_score(
  true_values,
  lower,
  upper,
  interval_range,
  weigh = TRUE,
  separate_results = FALSE
)
```

**Arguments**

**true\_values** A vector with the true observed values of size n

**lower** vector of size n with the prediction for the lower quantile of the given range

**upper** vector of size n with the prediction for the upper quantile of the given range

**interval\_range** the range of the prediction intervals. i.e. if you're forecasting the 0.05 and 0.95 quantile, the interval\_range would be 90. Can be either a single number or a vector of size n, if the range changes for different forecasts to be scored. This corresponds to  $(100-\alpha)/100$  in Gneiting and Raftery (2007). Internally, the range will be transformed to alpha.

**weigh** if TRUE, weigh the score by  $\alpha / 2$ , so it can be averaged into an interval score that, in the limit, corresponds to CRPS. Alpha is the decimal value that represents how much is outside a central prediction interval (e.g. for a 90 percent central prediction interval, alpha is 0.1) Default: TRUE.

**separate\_results** if TRUE (default is FALSE), then the separate parts of the interval score (dispersion penalty, penalties for over- and under-prediction get returned as separate elements of a list). If you want a `data.frame` instead, simply call `as.data.frame()` on the output.

**Value**

vector with the scoring values, or a list with separate entries if `separate_results` is TRUE.

**References**

Strictly Proper Scoring Rules, Prediction, and Estimation, Tilmann Gneiting and Adrian E. Raftery, 2007, Journal of the American Statistical Association, Volume 102, 2007 - Issue 477

Evaluating epidemic forecasts in an interval format, Johannes Bracher, Evan L. Ray, Tilmann Gneiting and Nicholas G. Reich, <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1008618> # nolint

**Examples**

```
true_values <- rnorm(30, mean = 1:30)
interval_range <- rep(90, 30)
alpha <- (100 - interval_range) / 100
```



```
lower <- qnorm(alpha / 2, rnorm(30, mean = 1:30))
upper <- qnorm((1 - alpha / 2), rnorm(30, mean = 1:30))

interval_score(
  true_values = true_values,
  lower = lower,
  upper = upper,
  interval_range = interval_range
)

# gives a warning, as the interval_range should likely be 50 instead of 0.5
interval_score(true_value = 4, upper = 2, lower = 8, interval_range = 0.5)

# example with missing values and separate results
interval_score(
  true_values = c(true_values, NA),
  lower = c(lower, NA),
  upper = c(NA, upper),
  separate_results = TRUE,
  interval_range = 90
)
```

---

logs\_binary

*Log Score for Binary outcomes*

---

### Description

Computes the Log Score for probabilistic forecasts of binary outcomes.

### Usage

```
logs_binary(true_values, predictions)
```

### Arguments

true_values	A vector with the true observed values of size n with all values equal to either 0 or 1
predictions	A vector with a predicted probability that true_value = 1.

### Details

The Log Score is a proper score rule suited to assessing the accuracy of probabilistic binary predictions. The outcomes can be either 0 or 1, the predictions must be a probability that the true outcome will be 1.

The Log Score is then computed as the negative logarithm of the probability assigned to the true outcome. Reporting the negative logarithm means that smaller values are better.

**Value**

A numeric value with the Log Score, i.e. the mean squared error of the given probability forecasts

**Examples**

```
true_values <- sample(c(0, 1), size = 30, replace = TRUE)
predictions <- runif(n = 30, min = 0, max = 1)
logs_binary(true_values, predictions)
```

---

logs\_sample

*Logarithmic score*

---

**Description**

Wrapper around the `logs_sample()` function from the **scoringRules** package. Used to score continuous predictions. While the Log Score is in theory also applicable to integer forecasts, the problem lies in the implementation: The Log Score needs a kernel density estimation, which is not well defined with integer-valued Monte Carlo Samples. The Log Score can be used for specific integer valued probability distributions. See the scoringRules package for more details.

**Usage**

```
logs_sample(true_values, predictions)
```

**Arguments**

true_values	A vector with the true observed values of size n
predictions	nxN matrix of predictive samples, n (number of rows) being the number of data points and N (number of columns) the number of Monte Carlo samples. Alternatively, predictions can just be a vector of size n.

**Value**

vector with the scoring values

**References**

Alexander Jordan, Fabian Krüger, Sebastian Lerch, Evaluating Probabilistic Forecasts with scoringRules, <https://www.jstatsoft.org/article/view/v090i12>

**Examples**

```
true_values <- rpois(30, lambda = 1:30)
predictions <- replicate(200, rpois(n = 30, lambda = 1:30))
logs_sample(true_values, predictions)
```

---

`log_shift`*Log transformation with an additive shift*

---

**Description**

Function that shifts a value by some offset and then applies the natural logarithm to it.

**Usage**

```
log_shift(x, offset = 0, base = exp(1))
```

**Arguments**

<code>x</code>	vector of input values to be transformed
<code>offset</code>	number to add to the input value before taking the natural logarithm
<code>base</code>	a positive or complex number: the base with respect to which logarithms are computed. Defaults to $e = \exp(1)$ .

**Details**

The output is computed as  $\log(x + \text{offset})$

**Value**

A numeric vector with transformed values

**References**

Transformation of forecasts for evaluating predictive performance in an epidemiological context  
Nikos I. Bosse, Sam Abbott, Anne Cori, Edwin van Leeuwen, Johannes Bracher, Sebastian Funk  
medRxiv 2023.01.23.23284722 doi:10.1101/2023.01.23.23284722 <https://www.medrxiv.org/content/10.1101/2023.01.23.23284722v1> # nolint

**Examples**

```
log_shift(1:10)
log_shift(0:9, offset = 1)

transform_forecasts(
  example_quantile[true_value > 0, ],
  fun = log_shift,
  offset = 1
)
```

---

`mad_sample`*Determine dispersion of a probabilistic forecast*

---

## Description

Determine dispersion of a probabilistic forecast

## Usage

```
mad_sample(predictions)
```

## Arguments

`predictions`      nxN matrix of predictive samples, n (number of rows) being the number of data points and N (number of columns) the number of Monte Carlo samples. Alternatively, predictions can just be a vector of size n.

## Details

Sharpness is the ability of the model to generate predictions within a narrow range and dispersion is the lack thereof. It is a data-independent measure, and is purely a feature of the forecasts themselves.

Dispersion of predictive samples corresponding to one single true value is measured as the normalised median of the absolute deviation from the median of the predictive samples. For details, see [mad\(\)](#) and the explanations given in Funk et al. (2019)

## Value

vector with dispersion values

## References

Funk S, Camacho A, Kucharski AJ, Lowe R, Eggo RM, Edmunds WJ (2019) Assessing the performance of real-time epidemic forecasts: A case study of Ebola in the Western Area region of Sierra Leone, 2014-15. PLoS Comput Biol 15(2): e1006785. [doi:10.1371/journal.pcbi.1006785](https://doi.org/10.1371/journal.pcbi.1006785)

## Examples

```
predictions <- replicate(200, rpois(n = 30, lambda = 1:30))
mad_sample(predictions)
```

make\_NA

*Make Rows NA in Data for Plotting***Description**

Filters the data and turns values into NA before the data gets passed to `plot_predictions()`. The reason to do this is to this is that it allows to 'filter' prediction and truth data separately. Any value that is NA will then be removed in the subsequent call to `plot_predictions()`.

**Usage**

```
make_NA(data = NULL, what = c("truth", "forecast", "both"), ...)
```

```
make_na(data = NULL, what = c("truth", "forecast", "both"), ...)
```

**Arguments**

**data** A data.frame or data.table with the predictions and observations. For scoring using `score()`, the following columns need to be present:

- `true_value` - the true observed values
- `prediction` - predictions or predictive samples for one true value. (You only don't need to provide a prediction column if you want to score quantile forecasts in a wide range format.)

For scoring integer and continuous forecasts a `sample` column is needed:

- `sample` - an index to identify the predictive samples in the prediction column generated by one model for one true value. Only necessary for continuous and integer forecasts, not for binary predictions.

For scoring predictions in a quantile-format forecast you should provide a column called `quantile`:

- `quantile`: quantile to which the prediction corresponds

In addition a `model` column is suggested and if not present this will be flagged and added to the input data with all forecasts assigned as an "unspecified model").

You can check the format of your data using `check_forecasts()` and there are examples for each format (`example_quantile`, `example_continuous`, `example_integer`, and `example_binary`).

**what** character vector that determines which values should be turned into NA. If `what = "truth"`, values in the column 'true\_value' will be turned into NA. If `what = "forecast"`, values in the column 'prediction' will be turned into NA. If `what = "both"`, values in both column will be turned into NA.

**...** logical statements used to filter the data

**Value**

A data.table

## Examples

```
make_NA (  
  example_continuous,  
  what = "truth",  
  target_end_date >= "2021-07-22",  
  target_end_date < "2021-05-01"  
)
```

---

merge\_pred\_and\_obs      *Merge Forecast Data And Observations*

---

## Description

The function more or less provides a wrapper around `merge` that aims to handle the merging well if additional columns are present in one or both data sets. If in doubt, you should probably merge the data sets manually.

## Usage

```
merge_pred_and_obs(  
  forecasts,  
  observations,  
  join = c("left", "full", "right"),  
  by = NULL  
)
```

## Arguments

<code>forecasts</code>	data.frame with the forecast data (as can be passed to <code>score()</code> ).
<code>observations</code>	data.frame with the observations
<code>join</code>	character, one of <code>c("left", "full", "right")</code> . Determines the type of the join. Usually, a left join is appropriate, but sometimes you may want to do a full join to keep dates for which there is a forecast, but no ground truth data.
<code>by</code>	character vector that denotes the columns by which to merge. Any value that is not a column in observations will be removed.

## Value

a data.frame with forecasts and observations

## Examples

```
forecasts <- example_quantile_forecasts_only  
observations <- example_truth_only  
merge_pred_and_obs(forecasts, observations)
```

---

`metrics`*Summary information for selected metrics*

---

### Description

A data set with summary information on selected metrics implemented in **scoringutils**

### Usage

```
metrics
```

### Format

An object of class `data.table` (inherits from `data.frame`) with 22 rows and 8 columns.

### Details

The data was created using the script `create-metric-tables.R` in the `inst/` folder (or the top level folder in a compiled package).

---

`pairwise_comparison`*Do Pairwise Comparisons of Scores*

---

### Description

Compute relative scores between different models making pairwise comparisons. Pairwise comparisons are a sort of pairwise tournament where all combinations of two models are compared against each other based on the overlapping set of available forecasts common to both models. Internally, a ratio of the mean scores of both models is computed. The relative score of a model is then the geometric mean of all mean score ratios which involve that model. When a baseline is provided, then that baseline is excluded from the relative scores for individual models (which therefore differ slightly from relative scores without a baseline) and all relative scores are scaled by (i.e. divided by) the relative score of the baseline model. Usually, the function input should be unsummarised scores as produced by `score()`. Note that the function internally infers the *unit of a single forecast* by determining all columns in the input that do not correspond to metrics computed by `score()`. Adding unrelated columns will change results in an unpredictable way.

The code for the pairwise comparisons is inspired by an implementation by Johannes Bracher. The implementation of the permutation test follows the function `permutationTest` from the `surveillance` package by Michael Höhle, Andrea Riebler and Michaela Paul.

**Usage**

```
pairwise_comparison(
  scores,
  by = "model",
  metric = "auto",
  baseline = NULL,
  ...
)
```

**Arguments**

scores	A data.table of scores as produced by <a href="#">score()</a> .
by	character vector with names of columns present in the input data.frame. by determines how pairwise comparisons will be computed. You will get a relative skill score for every grouping level determined in by. If, for example, by = c("model", "location"). Then you will get a separate relative skill score for every model in every location. Internally, the data.frame will be split according by (but removing "model" before splitting) and the pairwise comparisons will be computed separately for the split data.frames.
metric	A character vector of length one with the metric to do the comparison on. The default is "auto", meaning that either "interval_score", "crps", or "brier_score" will be selected where available. See <a href="#">available_metrics()</a> for available metrics.
baseline	character vector of length one that denotes the baseline model against which to compare other models.
...	additional arguments for the comparison between two models. See <a href="#">compare_two_models()</a> for more information.

**Value**

A ggplot2 object with a coloured table of summarised scores

**Author(s)**

Nikos Bosse <nikosbosse@gmail.com>  
 Johannes Bracher, <johannes.bracher@kit.edu>

**Examples**

```
scores <- score(example_quantile)
pairwise <- pairwise_comparison(scores, by = "target_type")

library(ggplot2)
plot_pairwise_comparison(pairwise, type = "mean_scores_ratio") +
  facet_wrap(~target_type)
```



---

pit

*Probability Integral Transformation (data.frame Format)*

---

## Description

Wrapper around `pit()` for use in data.frames

## Usage

```
pit(data, by, n_replicates = 100)
```

## Arguments

`data` a data.frame with the following columns: `true_value`, `prediction`, `sample`.

`by` Character vector with the columns according to which the PIT values shall be grouped. If you e.g. have the columns 'model' and 'location' in the data and want to have a PIT histogram for every model and location, specify `by = c("model", "location")`.

`n_replicates` the number of draws for the randomised PIT for integer predictions.

## Details

see [pit\(\)](#)

## Value

a data.table with PIT values according to the grouping specified in `by`

## References

Sebastian Funk, Anton Camacho, Adam J. Kucharski, Rachel Lowe, Rosalind M. Eggo, W. John Edmunds (2019) Assessing the performance of real-time epidemic forecasts: A case study of Ebola in the Western Area region of Sierra Leone, 2014-15, [doi:10.1371/journal.pcbi.1006785](https://doi.org/10.1371/journal.pcbi.1006785)

## Examples

```
result <- pit(example_continuous, by = "model")
plot_pit(result)

# example with quantile data
result <- pit(example_quantile, by = "model")
plot_pit(result)
```

---

 pit\_sample

*Probability Integral Transformation (sample-based version)*


---

### Description

Uses a Probability Integral Transformation (PIT) (or a randomised PIT for integer forecasts) to assess the calibration of predictive Monte Carlo samples. Returns a p-values resulting from an Anderson-Darling test for uniformity of the (randomised) PIT as well as a PIT histogram if specified.

### Usage

```
pit_sample(true_values, predictions, n_replicates = 100)
```

### Arguments

true_values	A vector with the true observed values of size n
predictions	nxN matrix of predictive samples, n (number of rows) being the number of data points and N (number of columns) the number of Monte Carlo samples. Alternatively, predictions can just be a vector of size n.
n_replicates	the number of draws for the randomised PIT for integer predictions.

### Details

Calibration or reliability of forecasts is the ability of a model to correctly identify its own uncertainty in making predictions. In a model with perfect calibration, the observed data at each time point look as if they came from the predictive probability distribution at that time.

Equivalently, one can inspect the probability integral transform of the predictive distribution at time  $t$ ,

$$u_t = F_t(x_t)$$

where  $x_t$  is the observed data point at time  $t$  in  $t_1, \dots, t_n$ ,  $n$  being the number of forecasts, and  $F_t$  is the (continuous) predictive cumulative probability distribution at time  $t$ . If the true probability distribution of outcomes at time  $t$  is  $G_t$  then the forecasts  $F_t$  are said to be ideal if  $F_t = G_t$  at all times  $t$ . In that case, the probabilities  $u_t$  are distributed uniformly.

In the case of discrete outcomes such as incidence counts, the PIT is no longer uniform even when forecasts are ideal. In that case a randomised PIT can be used instead:

$$u_t = P_t(k_t) + v * (P_t(k_t) - P_t(k_t - 1))$$

where  $k_t$  is the observed count,  $P_t(x)$  is the predictive cumulative probability of observing incidence  $k$  at time  $t$ ,  $P_t(-1) = 0$  by definition and  $v$  is standard uniform and independent of  $k$ . If  $P_t$  is the true cumulative probability distribution, then  $u_t$  is standard uniform.

The function checks whether integer or continuous forecasts were provided. It then applies the (randomised) probability integral and tests the values  $u_t$  for uniformity using the Anderson-Darling test.

As a rule of thumb, there is no evidence to suggest a forecasting model is miscalibrated if the p-value found was greater than a threshold of  $p \geq 0.1$ , some evidence that it was miscalibrated if  $0.01 < p < 0.1$ , and good evidence that it was miscalibrated if  $p \leq 0.01$ . However, the AD-p-values may be overly strict and their actual usefulness may be questionable. In this context it should be noted, though, that uniformity of the PIT is a necessary but not sufficient condition of calibration.

### Value

A vector with PIT-values. For continuous forecasts, the vector will correspond to the length of `true_values`. For integer forecasts, a randomised PIT will be returned of length `length(true_values) * n_replicates`

### References

Sebastian Funk, Anton Camacho, Adam J. Kucharski, Rachel Lowe, Rosalind M. Eggo, W. John Edmunds (2019) Assessing the performance of real-time epidemic forecasts: A case study of Ebola in the Western Area region of Sierra Leone, 2014-15, [doi:10.1371/journal.pcbi.1006785](https://doi.org/10.1371/journal.pcbi.1006785)

### See Also

[pit\(\)](#)

### Examples

```
## continuous predictions
true_values <- rnorm(20, mean = 1:20)
predictions <- replicate(100, rnorm(n = 20, mean = 1:20))
pit <- pit_sample(true_values, predictions)
plot_pit(pit)

## integer predictions
true_values <- rpois(50, lambda = 1:50)
predictions <- replicate(2000, rpois(n = 50, lambda = 1:50))
pit <- pit_sample(true_values, predictions, n_replicates = 30)
plot_pit(pit)
```

---

plot\_avail\_forecasts *Visualise Where Forecasts Are Available*

---

### Description

Visualise Where Forecasts Are Available

**Usage**

```
plot_avail_forecasts(
  avail_forecasts,
  y = "model",
  x = "forecast_date",
  make_x_factor = TRUE,
  show_numbers = TRUE
)
```

**Arguments**

avail_forecasts	data.frame with a column called Number forecasts as produced by <a href="#">avail_forecasts()</a>
y	character vector of length one that denotes the name of the column to appear on the y-axis of the plot. Default is "model".
x	character vector of length one that denotes the name of the column to appear on the x-axis of the plot. Default is "forecast_date".
make_x_factor	logical (default is TRUE). Whether or not to convert the variable on the x-axis to a factor. This has an effect e.g. if dates are shown on the x-axis.
show_numbers	logical (default is TRUE) that indicates whether or not to show the actual count numbers on the plot

**Value**

ggplot object with a plot of interval coverage

**Examples**

```
library(ggplot2)
avail_forecasts <- avail_forecasts(
  example_quantile, by = c("model", "target_type", "target_end_date")
)
plot_avail_forecasts(
  avail_forecasts, x = "target_end_date", show_numbers = FALSE
) +
  facet_wrap("target_type")
```

---

plot\_correlation      *Plot Correlation Between Metrics*

---

**Description**

Plots a heatmap of correlations between different metrics

**Usage**

```
plot_correlation(correlations)
```

**Arguments**

correlations    A data.table of correlations between scores as produced by `correlation()`.

**Value**

A ggplot2 object showing a coloured matrix of correlations between metrics

**Examples**

```
scores <- score(example_quantile)
correlations <- correlation(
  summarise_scores(scores)
)
plot_correlation(correlations)
```

---

plot\_heatmap

*Create a Heatmap of a Scoring Metric*

---

**Description**

This function can be used to create a heatmap of one metric across different groups, e.g. the interval score obtained by several forecasting models in different locations.

**Usage**

```
plot_heatmap(scores, y = "model", x, metric)
```

**Arguments**

scores            A data.frame of scores based on quantile forecasts as produced by `score()`.

y                 The variable from the scores you want to show on the y-Axis. The default for this is "model"

x                 The variable from the scores you want to show on the x-Axis. This could be something like "horizon", or "location"

metric            the metric that determines the value and colour shown in the tiles of the heatmap

**Value**

A ggplot2 object showing a heatmap of the desired metric

**Examples**

```
scores <- score(example_quantile)
scores <- summarise_scores(scores, by = c("model", "target_type", "range"))

plot_heatmap(scores, x = "target_type", metric = "bias")
```

---

plot\_interval\_coverage

*Plot Interval Coverage*

---

### Description

Plot interval coverage

### Usage

```
plot_interval_coverage(scores, colour = "model")
```

### Arguments

scores	A data.frame of scores based on quantile forecasts as produced by <code>score()</code> or <code>summarise_scores()</code> . Note that "range" must be included in the by argument when running <code>summarise_scores()</code>
colour	According to which variable shall the graphs be coloured? Default is "model".

### Value

ggplot object with a plot of interval coverage

### Examples

```
scores <- score(example_quantile)
scores <- summarise_scores(scores, by = c("model", "range"))
plot_interval_coverage(scores)
```

---

plot\_pairwise\_comparison

*Plot Heatmap of Pairwise Comparisons*

---

### Description

Creates a heatmap of the ratios or pvalues from a pairwise comparison between models

### Usage

```
plot_pairwise_comparison(
  comparison_result,
  type = c("mean_scores_ratio", "pval")
)
```

**Arguments**

comparison_result	A data.frame as produced by <code>pairwise_comparison()</code>
type	character vector of length one that is either "mean_scores_ratio" or "pval". This denotes whether to visualise the ratio or the p-value of the pairwise comparison. Default is "mean_scores_ratio".

**Examples**

```
library(ggplot2)
scores <- score(example_quantile)
pairwise <- pairwise_comparison(scores, by = "target_type")
plot_pairwise_comparison(pairwise, type = "mean_scores_ratio") +
  facet_wrap(~target_type)
```

---

plot\_pit

*PIT Histogram*


---

**Description**

Make a simple histogram of the probability integral transformed values to visually check whether a uniform distribution seems likely.

**Usage**

```
plot_pit(pit, num_bins = "auto", breaks = NULL)
```

**Arguments**

pit	either a vector with the PIT values of size n, or a data.frame as produced by <code>pit()</code>
num_bins	the number of bins in the PIT histogram, default is "auto". When num_bins == "auto", <code>plot_pit()</code> will either display 10 bins, or it will display a bin for each available quantile in case you passed in data in a quantile-based format. You can control the number of bins by supplying a number. This is fine for sample-based pit histograms, but may fail for quantile-based formats. In this case it is preferred to supply explicit breaks points using the breaks argument.
breaks	numeric vector with the break points for the bins in the PIT histogram. This is preferred when creating a PIT histogram based on quantile-based data. Default is NULL and breaks will be determined by num_bins.

**Value**

vector with the scoring values

## Examples

```
# PIT histogram in vector based format
true_values <- rnorm(30, mean = 1:30)
predictions <- replicate(200, rnorm(n = 30, mean = 1:30))
pit <- pit_sample(true_values, predictions)
plot_pit(pit)

# quantile-based pit
pit <- pit(example_quantile, by = "model")
plot_pit(pit, breaks = seq(0.1, 1, 0.1))

# sample-based pit
pit <- pit(example_integer, by = "model")
plot_pit(pit)
```

---

plot\_predictions      *Plot Predictions vs True Values*

---

## Description

Make a plot of observed and predicted values

## Usage

```
plot_predictions(data, by = NULL, x = "date", range = c(0, 50, 90))
```

## Arguments

data	a data.frame that follows the same specifications outlined in <a href="#">score()</a> . To customise your plotting, you can filter your data using the function <a href="#">make_NA()</a> .
by	character vector with column names that denote categories by which the plot should be stratified. If for example you want to have a faceted plot, this should be a character vector with the columns used in facetting (note that the facetting still needs to be done outside of the function call)
x	character vector of length one that denotes the name of the variable
range	numeric vector indicating the interval ranges to plot. If 0 is included in range, the median prediction will be shown.

## Value

ggplot object with a plot of true vs predicted values



**Examples**

```
library(ggplot2)
library(magrittr)

example_continuous %>%
  make_NA (
    what = "truth",
    target_end_date >= "2021-07-22",
    target_end_date < "2021-05-01"
  ) %>%
  make_NA (
    what = "forecast",
    model != "EuroCOVIDhub-ensemble",
    forecast_date != "2021-06-07"
  ) %>%
  plot_predictions (
    x = "target_end_date",
    by = c("target_type", "location"),
    range = c(0, 50, 90, 95)
  ) +
  facet_wrap(~ location + target_type, scales = "free_y") +
  aes(fill = model, color = model)

example_continuous %>%
  make_NA (
    what = "truth",
    target_end_date >= "2021-07-22",
    target_end_date < "2021-05-01"
  ) %>%
  make_NA (
    what = "forecast",
    forecast_date != "2021-06-07"
  ) %>%
  plot_predictions (
    x = "target_end_date",
    by = c("target_type", "location"),
    range = c(0)
  ) +
  facet_wrap(~ location + target_type, scales = "free_y") +
  aes(fill = model, color = model)
```

---

plot\_quantile\_coverage

*Plot Quantile Coverage*

---

**Description**

Plot quantile coverage

**Usage**

```
plot_quantile_coverage(scores, colour = "model")
```

**Arguments**

**scores** A data.frame of scores based on quantile forecasts as produced by `score()` or `summarise_scores()`. Note that "range" must be included in the by argument when running `summarise_scores()`

**colour** According to which variable shall the graphs be coloured? Default is "model".

**Value**

ggplot object with a plot of interval coverage

**Examples**

```
scores <- score(example_quantile)
scores <- summarise_scores(scores, by = c("model", "quantile"))
plot_quantile_coverage(scores)
```

---

plot\_ranges

*Plot Metrics by Range of the Prediction Interval*

---

**Description**

Visualise the metrics by range, e.g. if you are interested how different interval ranges contribute to the overall interval score, or how sharpness / dispersion changes by range.

**Usage**

```
plot_ranges(scores, y = "interval_score", x = "model", colour = "range")
```

**Arguments**

**scores** A data.frame of scores based on quantile forecasts as produced by `score()` or `summarise_scores()`. Note that "range" must be included in the by argument when running `summarise_scores()`

**y** The variable from the scores you want to show on the y-Axis. This could be something like "interval\_score" (the default) or "dispersion"

**x** The variable from the scores you want to show on the x-Axis. Usually this will be "model"

**colour** Character vector of length one used to determine a variable for colouring dots. The Default is "range".

**Value**

A ggplot2 object showing a contributions from the three components of the weighted interval score

**Examples**

```
library(ggplot2)
scores <- score(example_quantile)
scores <- summarise_scores(scores, by = c("model", "target_type", "range"))

plot_ranges(scores, x = "model") +
  facet_wrap(~target_type, scales = "free")

# visualise dispersion instead of interval score
plot_ranges(scores, y = "dispersion", x = "model") +
  facet_wrap(~target_type)
```

---

plot_score_table	<i>Plot Coloured Score Table</i>
------------------	----------------------------------

---

**Description**

Plots a coloured table of summarised scores obtained using `score()`.

**Usage**

```
plot_score_table(scores, y = "model", by = NULL, metrics = NULL)
```

**Arguments**

scores	A data.table of scores as produced by <code>score()</code> .
y	the variable to be shown on the y-axis. Instead of a single character string, you can also specify a vector with column names, e.g. <code>y = c("model", "location")</code> . These column names will be concatenated to create a unique row identifier (e.g. "model1_location1").
by	A character vector that determines how the colour shading for the plot gets computed. By default (NULL), shading will be determined per metric, but you can provide additional column names (see examples).
metrics	A character vector with the metrics to show. If set to NULL (default), all metrics present in scores will be shown.

**Value**

A ggplot2 object with a coloured table of summarised scores

**Examples**

```
library(ggplot2)
library(magrittr) # pipe operator

scores <- score(example_quantile) %>%
```

```

summarise_scores(by = c("model", "target_type")) %>%
summarise_scores(fun = signif, digits = 2)

plot_score_table(scores, y = "model", by = "target_type") +
  facet_wrap(~target_type, ncol = 1)

# can also put target description on the y-axis
plot_score_table(scores,
  y = c("model", "target_type"),
  by = "target_type")

```

---

plot\_wis

*Plot Contributions to the Weighted Interval Score*


---

### Description

Visualise the components of the weighted interval score: penalties for over-prediction, under-prediction and for high dispersion (lack of sharpness)

### Usage

```
plot_wis(scores, x = "model", relative_contributions = FALSE, flip = FALSE)
```

### Arguments

scores	A data.frame of scores based on quantile forecasts as produced by <code>score()</code> and summarised using <code>summarise_scores()</code>
x	The variable from the scores you want to show on the x-Axis. Usually this will be "model".
relative_contributions	show relative contributions instead of absolute contributions. Default is FALSE and this functionality is not available yet.
flip	boolean (default is FALSE), whether or not to flip the axes.

### Value

A ggplot2 object showing a contributions from the three components of the weighted interval score

### References

Bracher J, Ray E, Gneiting T, Reich, N (2020) Evaluating epidemic forecasts in an interval format. <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1008618>

**Examples**

```
library(ggplot2)
scores <- score(example_quantile)
scores <- summarise_scores(scores, by = c("model", "target_type"))

plot_wis(scores,
  x = "model",
  relative_contributions = TRUE
) +
  facet_wrap(~target_type)
plot_wis(scores,
  x = "model",
  relative_contributions = FALSE
) +
  facet_wrap(~target_type, scales = "free_x")
```

---

```
print.scoringutils_check
```

*Print output from check\_forecasts()*

---

**Description**

Helper function that prints the output generated by [check\\_forecasts\(\)](#)

**Usage**

```
## S3 method for class 'scoringutils_check'
print(x, ...)
```

**Arguments**

x	An object of class 'scoringutils_check' as produced by <a href="#">check_forecasts()</a>
...	additional arguments (not used here)

**Examples**

```
check <- check_forecasts(example_quantile)
print(check)
```

---

quantile_score	<i>Quantile Score</i>
----------------	-----------------------

---

### Description

Proper Scoring Rule to score quantile predictions. Smaller values are better. The quantile score is closely related to the Interval score (see `interval_score()`) and is the quantile equivalent that works with single quantiles instead of central prediction intervals.

### Usage

```
quantile_score(true_values, predictions, quantiles, weigh = TRUE)
```

### Arguments

<code>true_values</code>	A vector with the true observed values of size <code>n</code>
<code>predictions</code>	<code>n</code> × <code>N</code> matrix of predictive samples, <code>n</code> (number of rows) being the number of data points and <code>N</code> (number of columns) the number of Monte Carlo samples. Alternatively, predictions can just be a vector of size <code>n</code> .
<code>quantiles</code>	vector of size <code>n</code> with the quantile values of the corresponding predictions.
<code>weigh</code>	if TRUE, weigh the score by $\alpha / 2$ , so it can be averaged into an interval score that, in the limit, corresponds to CRPS. Alpha is the value that corresponds to the $(\alpha/2)$ or $(1 - \alpha/2)$ quantiles provided and will be computed from the quantile. Alpha is the decimal value that represents how much is outside a central prediction interval (E.g. for a 90 percent central prediction interval, alpha is 0.1). Default: TRUE.

### Value

vector with the scoring values

### References

Strictly Proper Scoring Rules, Prediction, and Estimation, Tilmann Gneiting and Adrian E. Raftery, 2007, Journal of the American Statistical Association, Volume 102, 2007 - Issue 477

Evaluating epidemic forecasts in an interval format, Johannes Bracher, Evan L. Ray, Tilmann Gneiting and Nicholas G. Reich, <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1008618>

### Examples

```
true_values <- rnorm(10, mean = 1:10)
alpha <- 0.5

lower <- qnorm(alpha / 2, rnorm(10, mean = 1:10))
upper <- qnorm((1 - alpha / 2), rnorm(10, mean = 1:10))
```

```
qs_lower <- quantile_score(true_values,  
  predictions = lower,  
  quantiles = alpha / 2  
)  
qs_upper <- quantile_score(true_values,  
  predictions = upper,  
  quantiles = 1 - alpha / 2  
)  
interval_score <- (qs_lower + qs_upper) / 2
```

---

sample\_to\_quantile      *Change Data from a Sample Based Format to a Quantile Format*

---

## Description

Transform data from a format that is based on predictive samples to a format based on plain quantiles.

## Usage

```
sample_to_quantile(data, quantiles = c(0.05, 0.25, 0.5, 0.75, 0.95), type = 7)
```

## Arguments

**data**                    a data.frame with samples

**quantiles**                a numeric vector of quantiles to extract

**type**                    type argument passed down to the quantile function. For more information, see [quantile\(\)](#)

## Value

a data.frame in a long interval range format

## Examples

```
sample_to_quantile(example_integer)
```

score

*Evaluate forecasts*

## Description

This function allows automatic scoring of forecasts using a range of metrics. For most users it will be the workhorse for scoring forecasts as it wraps the lower level functions package functions. However, these functions are also available if you wish to make use of them independently.

A range of forecasts formats are supported, including quantile-based, sample-based, binary forecasts. Prior to scoring, users may wish to make use of `check_forecasts()` to ensure that the input data is in a supported format though this will also be run internally by `score()`. Examples for each format are also provided (see the documentation for data below or in `check_forecasts()`).

Each format has a set of required columns (see below). Additional columns may be present to indicate a grouping of forecasts. For example, we could have forecasts made by different models in various locations at different time points, each for several weeks into the future. It is important, that there are only columns present which are relevant in order to group forecasts. A combination of different columns should uniquely define the *unit of a single forecast*, meaning that a single forecast is defined by the values in the other columns. Adding additional unrelated columns may alter results.

To obtain a quick overview of the currently supported evaluation metrics, have a look at the [metrics](#) data included in the package. The column `metrics$Name` gives an overview of all available metric names that can be computed. If interested in an unsupported metric please open a [feature request](#) or consider contributing a pull request.

For additional help and examples, check out the [Getting Started Vignette](#) as well as the paper [Evaluating Forecasts with scoringutils in R](#).

## Usage

```
score(data, metrics = NULL, ...)
```

## Arguments

**data** A `data.frame` or `data.table` with the predictions and observations. For scoring using `score()`, the following columns need to be present:

- `true_value` - the true observed values
- `prediction` - predictions or predictive samples for one true value. (You only don't need to provide a prediction column if you want to score quantile forecasts in a wide range format.)

For scoring integer and continuous forecasts a `sample` column is needed:

- `sample` - an index to identify the predictive samples in the prediction column generated by one model for one true value. Only necessary for continuous and integer forecasts, not for binary predictions.

For scoring predictions in a quantile-format forecast you should provide a column called `quantile`:



- `quantile`: quantile to which the prediction corresponds

In addition a `model` column is suggested and if not present this will be flagged and added to the input data with all forecasts assigned as an "unspecified model"). You can check the format of your data using `check_forecasts()` and there are examples for each format (`example_quantile`, `example_continuous`, `example_integer`, and `example_binary`).

`metrics` the metrics you want to have in the output. If NULL (the default), all available metrics will be computed. For a list of available metrics see `available_metrics()`, or check the `metrics` data set.

`...` additional parameters passed down to `score_quantile()` (internal function used for scoring forecasts in a quantile-based format).

### Value

A `data.table` with unsummarised scores. There will be one score per quantile or sample, which is usually not desired, so you should almost always run `summarise_scores()` on the unsummarised scores.

### Author(s)

Nikos Bosse <nikosbosse@gmail.com>

### References

Funk S, Camacho A, Kucharski AJ, Lowe R, Eggo RM, Edmunds WJ (2019) Assessing the performance of real-time epidemic forecasts: A case study of Ebola in the Western Area region of Sierra Leone, 2014-15. *PLoS Comput Biol* 15(2): e1006785. doi:10.1371/journal.pcbi.1006785

### Examples

```
library(magrittr) # pipe operator

check_forecasts(example_quantile)
score(example_quantile) %>%
  add_coverage(by = c("model", "target_type")) %>%
  summarise_scores(by = c("model", "target_type"))

# set forecast unit manually (to avoid issues with scoringutils trying to
# determine the forecast unit automatically), check forecasts before scoring
example_quantile %>%
  set_forecast_unit(
    c("location", "target_end_date", "target_type", "horizon", "model")
  ) %>%
  check_forecasts() %>%
  score()

# forecast formats with different metrics
## Not run:
score(example_binary)
```

```

score(example_quantile)
score(example_integer)
score(example_continuous)

## End(Not run)

# score point forecasts (marked by 'NA' in the quantile column)
score(example_point) %>%
  summarise_scores(by = "model", na.rm = TRUE)

```

---

set\_forecast\_unit      *Set unit of a single forecast manually*

---

## Description

Helper function to set the unit of a single forecast (i.e. the combination of columns that uniquely define a single forecast) manually. This simple function keeps the columns specified in `forecast_unit` (plus additional protected columns, e.g. for true values, predictions or quantile levels) and removes duplicate rows. If not done manually, `scoringutils` attempts to determine the unit of a single forecast automatically by simply assuming that all column names are relevant to determine the forecast unit. This may lead to unexpected behaviour, so setting the forecast unit explicitly can help make the code easier to debug and easier to read. When used as part of a workflow, `set_forecast_unit()` can be directly piped into `check_forecasts()` to check everything is in order.

## Usage

```
set_forecast_unit(data, forecast_unit)
```

## Arguments

**data**                    A `data.frame` or `data.table` with the predictions and observations. For scoring using `score()`, the following columns need to be present:

- `true_value` - the true observed values
- `prediction` - predictions or predictive samples for one true value. (You only don't need to provide a prediction column if you want to score quantile forecasts in a wide range format.)

For scoring integer and continuous forecasts a `sample` column is needed:

- `sample` - an index to identify the predictive samples in the prediction column generated by one model for one true value. Only necessary for continuous and integer forecasts, not for binary predictions.

For scoring predictions in a quantile-format forecast you should provide a column called `quantile`:

- `quantile`: quantile to which the prediction corresponds

In addition a `model` column is suggested and if not present this will be flagged and added to the input data with all forecasts assigned as an "unspecified model"). You can check the format of your data using `check_forecasts()` and there are examples for each format ([example\\_quantile](#), [example\\_continuous](#), [example\\_integer](#), and [example\\_binary](#)).

`forecast_unit` Character vector with the names of the columns that uniquely identify a single forecast.

### Value

A `data.table` with only those columns kept that are relevant to scoring or denote the unit of a single forecast as specified by the user.

### Examples

```
set_forecast_unit(
  example_quantile,
  c("location", "target_end_date", "target_type", "horizon", "model")
)
```

---

<code>se_mean_sample</code>	<i>Squared Error of the Mean (Sample-based Version)</i>
-----------------------------	---

---

### Description

Squared error of the mean calculated as

$$\text{mean}(\text{true\_value} - \text{prediction})^2$$

### Usage

```
se_mean_sample(true_values, predictions)
```

### Arguments

`true_values` A vector with the true observed values of size `n`

`predictions` `n` × `N` matrix of predictive samples, `n` (number of rows) being the number of data points and `N` (number of columns) the number of Monte Carlo samples. Alternatively, `predictions` can just be a vector of size `n`.

### Value

vector with the scoring values

### See Also

[squared\\_error\(\)](#)

**Examples**

```

true_values <- rnorm(30, mean = 1:30)
predicted_values <- rnorm(30, mean = 1:30)
se_mean_sample(true_values, predicted_values)

```

---

squared_error	<i>Squared Error</i>
---------------	----------------------

---

**Description**

Squared Error SE calculated as

$$(\text{true\_values} - \text{predicted\_values})^2$$

**Usage**

```
squared_error(true_values, predictions)
```

**Arguments**

`true_values`     A vector with the true observed values of size n  
`predictions`     A vector with predicted values of size n

**Value**

vector with the scoring values

**Examples**

```

true_values <- rnorm(30, mean = 1:30)
predicted_values <- rnorm(30, mean = 1:30)
squared_error(true_values, predicted_values)

```

---

summarise_scores	<i>Summarise scores as produced by <a href="#">score()</a></i>
------------------	--

---

**Description**

Summarise scores as produced by [score\(\)](#)

**Usage**

```

summarise_scores(
  scores,
  by = NULL,
  across = NULL,
  fun = mean,
  relative_skill = FALSE,
  relative_skill_metric = "auto",
  metric = deprecated(),
  baseline = NULL,
  ...
)

summarize_scores(
  scores,
  by = NULL,
  across = NULL,
  fun = mean,
  relative_skill = FALSE,
  relative_skill_metric = "auto",
  metric = deprecated(),
  baseline = NULL,
  ...
)

```

**Arguments**

<code>scores</code>	A <code>data.table</code> of scores as produced by <code>score()</code> .
<code>by</code>	character vector with column names to summarise scores by. Default is <code>NULL</code> , meaning that the only summary that takes is place is summarising over samples or quantiles (in case of quantile-based forecasts), such that there is one score per forecast as defined by the <i>unit of a single forecast</i> (rather than one score for every sample or quantile). The <i>unit of a single forecast</i> is determined by the columns present in the input data that do not correspond to a metric produced by <code>score()</code> , which indicate indicate a grouping of forecasts (for example there may be one forecast per day, location and model). Adding additional, unrelated, columns may alter results in an unpredictable way.
<code>across</code>	character vector with column names from the vector of variables that define the <i>unit of a single forecast</i> (see above) to summarise scores across (meaning that the specified columns will be dropped). This is an alternative to specifying <code>by</code> directly. If <code>NULL</code> (default), then <code>by</code> will be used or inferred internally if also not specified. Only one of <code>across</code> and <code>by</code> may be used at a time.
<code>fun</code>	a function used for summarising scores. Default is <code>mean</code> .
<code>relative_skill</code>	logical, whether or not to compute relative performance between models based on pairwise comparisons. If <code>TRUE</code> (default is <code>FALSE</code> ), then a column called 'model' must be present in the input data. For more information on the computation of relative skill, see <code>pairwise_comparison()</code> . Relative skill will be

calculated for the aggregation level specified in by.

`relative_skill_metric` character with the name of the metric for which a relative skill shall be computed. If equal to 'auto' (the default), then this will be either interval score, CRPS or Brier score (depending on which of these is available in the input data)

`metric` **[Deprecated]** Deprecated in 1.1.0. Use `relative_skill_metric` instead.

`baseline` character string with the name of a model. If a baseline is given, then a scaled relative skill with respect to the baseline will be returned. By default (NULL), relative skill will not be scaled with respect to a baseline model.

`...` additional parameters that can be passed to the summary function provided to fun. For more information see the documentation of the respective function.

## Examples

```
library(magrittr) # pipe operator

scores <- score(example_continuous)
summarise_scores(scores)

# summarise over samples or quantiles to get one score per forecast
scores <- score(example_quantile)
summarise_scores(scores)

# get scores by model
summarise_scores(scores, by = "model")

# get scores by model and target type
summarise_scores(scores, by = c("model", "target_type"))

# Get scores summarised across horizon, forecast date, and target end date
summarise_scores(
  scores, across = c("horizon", "forecast_date", "target_end_date")
)

# get standard deviation
summarise_scores(scores, by = "model", fun = sd)

# round digits
summarise_scores(scores, by = "model") %>%
  summarise_scores(fun = signif, digits = 2)

# get quantiles of scores
# make sure to aggregate over ranges first
summarise_scores(scores,
  by = "model", fun = quantile,
  probs = c(0.25, 0.5, 0.75)
)

# get ranges
```

```
# summarise_scores(scores, by = "range")
```

---

```
theme_scoringutils    Scoringutils ggplot2 theme
```

---

### Description

A theme for ggplot2 plots used in scoringutils

### Usage

```
theme_scoringutils()
```

### Value

A ggplot2 theme

---

```
transform_forecasts    Transform forecasts and observed values
```

---

### Description

Function to transform forecasts and true values before scoring.

### Usage

```
transform_forecasts(data, fun = log_shift, append = TRUE, label = "log", ...)
```

### Arguments

**data** A data.frame or data.table with the predictions and observations. For scoring using `score()`, the following columns need to be present:

- `true_value` - the true observed values
- `prediction` - predictions or predictive samples for one true value. (You only don't need to provide a prediction column if you want to score quantile forecasts in a wide range format.)

For scoring integer and continuous forecasts a `sample` column is needed:

- `sample` - an index to identify the predictive samples in the prediction column generated by one model for one true value. Only necessary for continuous and integer forecasts, not for binary predictions.

For scoring predictions in a quantile-format forecast you should provide a column called `quantile`:

- `quantile`: quantile to which the prediction corresponds

In addition a `model` column is suggested and if not present this will be flagged and added to the input data with all forecasts assigned as an "unspecified model"). You can check the format of your data using `check_forecasts()` and there are examples for each format ([example\\_quantile](#), [example\\_continuous](#), [example\\_integer](#), and [example\\_binary](#)).

fun	A function used to transform both true values and predictions. The default function is <code>log_shift()</code> , a custom function that is essentially the same as <code>log()</code> , but has an additional arguments ( <code>offset</code> ) that allows you add an offset before applying the logarithm. This is often helpful as the natural log transformation is not defined at zero. A common, and pragmatic solution, is to add a small offset to the data before applying the log transformation. In our work application we have often used an offset of 1 but the precise value will depend on your application.
append	Logical, defaults to TRUE. Whether or not to append a transformed version of the data to the currently existing data (TRUE). If selected, the data gets transformed and appended to the existing data frame, making it possible to use the outcome directly in <code>score()</code> . An additional column, 'scale', gets created that denotes which rows or untransformed ('scale' has the value "natural") and which have been transformed ('scale' has the value passed to the argument <code>label</code> ).
label	A string for the newly created 'scale' column to denote the newly transformed values. Only relevant if <code>append = TRUE</code> .
...	Additional parameters to pass to the function you supplied. For the default option of <code>log_shift()</code> this could be the <code>offset</code> argument.

### Details

There are a few reasons, depending on the circumstances, for why this might be desirable (check out the linked reference for more info). In epidemiology, for example, it may be useful to log-transform incidence counts before evaluating forecasts using scores such as the weighted interval score (WIS) or the continuous ranked probability score (CRPS). Log-transforming forecasts and observations changes the interpretation of the score from a measure of absolute distance between forecast and observation to a score that evaluates a forecast of the exponential growth rate. Another motivation can be to apply a variance-stabilising transformation or to standardise incidence counts by population.

Note that if you want to apply a transformation, it is important to transform the forecasts and observations and then apply the score. Applying a transformation after the score risks losing propriety of the proper scoring rule.

### Value

A `data.table` with either a transformed version of the data, or one with both the untransformed and the transformed data. includes the original data as well as a transformation of the original data. There will be one additional column, 'scale', present which will be set to "natural" for the untransformed forecasts.

### Author(s)

Nikos Bosse <[nikosbosse@gmail.com](mailto:nikosbosse@gmail.com)>



## References

Transformation of forecasts for evaluating predictive performance in an epidemiological context  
Nikos I. Bosse, Sam Abbott, Anne Cori, Edwin van Leeuwen, Johannes Bracher, Sebastian Funk  
medRxiv 2023.01.23.23284722 doi:10.1101/2023.01.23.23284722 <https://www.medrxiv.org/content/10.1101/2023.01.23.23284722v1>

## Examples

```
library(magrittr) # pipe operator

# transform forecasts using the natural logarithm
# negative values need to be handled (here by replacing them with 0)
example_quantile %>%
  .[, true_value := ifelse(true_value < 0, 0, true_value)] %>%
  # Here we use the default function log_shift() which is essentially the same
  # as log(), but has an additional arguments (offset) that allows you add an
  # offset before applying the logarithm.
  transform_forecasts(append = FALSE) %>%
  head()

# alternatively, integrating the truncation in the transformation function:
example_quantile %>%
  transform_forecasts(
    fun = function(x) {log_shift(pmax(0, x))}, append = FALSE
  ) %>%
  head()

# specifying an offset for the log transformation removes the
# warning caused by zeros in the data
example_quantile %>%
  .[, true_value := ifelse(true_value < 0, 0, true_value)] %>%
  transform_forecasts(offset = 1, append = FALSE) %>%
  head()

# adding square root transformed forecasts to the original ones
example_quantile %>%
  .[, true_value := ifelse(true_value < 0, 0, true_value)] %>%
  transform_forecasts(fun = sqrt, label = "sqrt") %>%
  score() %>%
  summarise_scores(by = c("model", "scale"))

# adding multiple transformations
example_quantile %>%
  .[, true_value := ifelse(true_value < 0, 0, true_value)] %>%
  transform_forecasts(fun = log_shift, offset = 1) %>%
  transform_forecasts(fun = sqrt, label = "sqrt") %>%
  head()
```

# Index

- \* **check-forecasts**
  - avail\_forecasts, 7
  - check\_forecasts, 13
  - find\_duplicates, 23
  - log\_shift, 27
  - print.scoringutils\_check, 45
  - transform\_forecasts, 55
- \* **data-handling**
  - merge\_pred\_and\_obs, 30
  - sample\_to\_quantile, 47
  - set\_forecast\_unit, 50
- \* **datasets**
  - example\_binary, 17
  - example\_continuous, 18
  - example\_integer, 19
  - example\_point, 19
  - example\_quantile, 20
  - example\_quantile\_forecasts\_only, 21
  - example\_truth\_only, 22
- \* **info**
  - available\_metrics, 6
  - metrics, 31
- \* **metric**
  - abs\_error, 3
  - ae\_median\_quantile, 5
  - ae\_median\_sample, 6
  - bias\_quantile, 8
  - bias\_range, 9
  - bias\_sample, 11
  - brier\_score, 12
  - crps\_sample, 15
  - dss\_sample, 16
  - interval\_score, 23
  - logs\_binary, 25
  - logs\_sample, 26
  - mad\_sample, 28
  - pit\_sample, 34
  - quantile\_score, 46
  - se\_mean\_sample, 51
  - squared\_error, 52
- \* **plotting**
  - make\_NA, 29
  - theme\_scoringutils, 55
- \* **scoring**
  - add\_coverage, 4
  - correlation, 15
  - pairwise\_comparison, 31
  - pit, 33
  - summarise\_scores, 52
- abs\_error, 3
- abs\_error(), 5, 6
- add\_coverage, 4
- add\_coverage(), 4
- ae\_median\_quantile, 5
- ae\_median\_quantile(), 4, 6
- ae\_median\_sample, 6
- ae\_median\_sample(), 4, 5
- as.data.frame(), 24
- avail\_forecasts, 7
- avail\_forecasts(), 36
- available\_metrics, 6
- available\_metrics(), 32, 49
- bias\_quantile, 8
- bias\_range, 9
- bias\_sample, 11
- brier\_score, 12
- check\_forecasts, 13
- check\_forecasts(), 7, 14, 29, 45, 48, 49, 51, 56
- compare\_two\_models(), 32
- correlation, 15
- correlation(), 37
- crps\_sample, 15
- crps\_sample(), 15
- dss\_sample, 16

dss\_sample(), 16

example\_binary, 7, 13, 14, 17, 29, 49, 51, 56

example\_continuous, 7, 13, 14, 18, 29, 49, 51, 56

example\_integer, 7, 13, 14, 19, 29, 49, 51, 56

example\_point, 19

example\_quantile, 7, 13, 14, 20, 29, 49, 51, 56

example\_quantile\_forecasts\_only, 21

example\_truth\_only, 22

find\_duplicates, 23

interval\_score, 23

interval\_score(), 46

log(), 56

log\_shift, 27

log\_shift(), 56

logs\_binary, 25

logs\_sample, 26

logs\_sample(), 26

mad(), 28

mad\_sample, 28

make\_NA, 29

make\_na (make\_NA), 29

make\_NA(), 40

merge\_pred\_and\_obs, 30

metrics, 31, 48, 49

pairwise\_comparison, 31

pairwise\_comparison(), 39, 53

pit, 33

pit(), 33, 35, 39

pit\_sample, 34

plot\_avail\_forecasts, 35

plot\_correlation, 36

plot\_heatmap, 37

plot\_interval\_coverage, 38

plot\_pairwise\_comparison, 38

plot\_pit, 39

plot\_pit(), 39

plot\_predictions, 40

plot\_predictions(), 29

plot\_quantile\_coverage, 41

plot\_ranges, 42

plot\_score\_table, 43

plot\_wis, 44

print.scoringutils\_check, 45

quantile(), 47

quantile\_score, 46

quantile\_score(), 23

sample\_to\_quantile, 47

sample\_to\_quantile(), 14

score, 48

score(), 4, 5, 7, 13–15, 23, 29–32, 37, 38, 40, 42–44, 48, 50, 52, 53, 55, 56

score\_quantile(), 49

se\_mean\_sample, 51

set\_forecast\_unit, 50

squared\_error, 52

squared\_error(), 51

summarise\_scores, 52

summarise\_scores(), 38, 42, 44, 49

summarize\_scores (summarise\_scores), 52

theme\_scoringutils, 55

transform\_forecasts, 55