

# Package ‘ohenery’

May 9, 2026

**Maintainer** Steven E. Pav <shabbychef@gmail.com>

**Version** 0.1.4

**Date** 2025-09-07

**License** LGPL-3

**Title** Modeling of Ordinal Random Variables via Softmax Regression

**BugReports** <https://github.com/shabbychef/ohenery/issues>

**Description** Supports the modeling of ordinal random variables, like the outcomes of races, via Softmax regression, under the Harville <[doi:10.1080/01621459.1973.10482425](https://doi.org/10.1080/01621459.1973.10482425)> and Henery <[doi:10.1111/j.2517-6161.1981.tb01153.x](https://doi.org/10.1111/j.2517-6161.1981.tb01153.x)> models.

**Depends** R (>= 3.0.2)

**Imports** Rcpp (>= 0.12.3), maxLik, magrittr, methods, generics, dplyr

**LinkingTo** Rcpp

**Suggests** rlang, tidyr, forcats, microbenchmark, testthat, numDeriv, ggplot2, scales, knitr

**Encoding** UTF-8

**URL** <https://github.com/shabbychef/ohenery>

**Collate** 'broom.r' 'data.r' 'harsm.r' 'hensm.r' 'linodds.r' 'ohenery.r' 'RcppExports.R' 'regularization.r' 'rsm.r' 'utils.r'

**RoxygenNote** 7.3.2

**NeedsCompilation** yes

**Author** Steven E. Pav [aut, cre] (ORCID: <<https://orcid.org/0000-0002-4197-6195>>)

**Repository** CRAN

**Date/Publication** 2025-09-09 07:50:13 UTC

## Contents

as.linodds	2
best_picture	4
diving	6
erank	7
glance.harsm	9
harsm	10
harsmfit	14
harsm_invlinc	17
hensm	18
hensmfit	22
inv_smax	25
normalize	26
ohenery	27
ohenery-NEWS	28
race_data	29
rhenery	31
rsm	32
smax	35
smlik	36
tidy.harsm	39
<b>Index</b>	<b>41</b>

---

as.linodds	<i>An object for modeling linear odds.</i>
------------	--

---

### Description

A model for odds linear in some feature.

### Usage

```
as.linodds(object, formula, beta)
```

```
## S3 method for class 'linodds'
predict(
  object,
  newdata,
  type = c("eta", "mu", "erank"),
  na.action = na.pass,
  group = NULL,
  ...
)
```

```
## S3 method for class 'linodds'
coef(object, ...)
```

**Arguments**

object	some list-like object.
formula	an object of class " <code>formula</code> " (or one that can be coerced to that class): a symbolic description of the model to be fitted. The details of model specification are given under 'Details'.
beta	the fit coefficients.
newdata	a <code>data.frame</code> from which we can extract a model frame via the formula of the object.
type	indicates which prediction should be returned: eta The odds. mu The probability. erank The expected rank.
na.action	How to deal with missing values in <code>y</code> , <code>g</code> , <code>X</code> , <code>wt</code> , <code>eta0</code> .
group	the string name of the <code>group</code> variable in the data, or a bare character with the group name. The group indices need not be integers, but that is more efficient. They need not be sorted.
...	other arguments.

**Details**

An object which holds a formula, some fit coefficients  $\beta$  which fit in that formula to generate odds in odds space. The odds can then be converted, via `predict.linodds` to probabilities, or to expected ranks under the Harville model. Both [harsm](#) and [hensm](#) return objects of class `linodds`.

We think of linear odds as  $\eta = x^T \beta$ , for independent variables  $x$ . The odds,  $\eta$  are converted to probabilities,  $\mu$  via  $\mu = c \exp \eta$ , where the constant  $c$  is chosen so the  $\mu$  for a given matching sum to one.

**Author(s)**

Steven E. Pav <[shabbychef@gmail.com](mailto:shabbychef@gmail.com)>

**See Also**

[harsm](#), [hensm](#).

[smax](#), [harsm\\_invlink](#).

---

 best\_picture

*Oscar Award Best Picture Data*


---

### Description

Historical data on the Best Picture nominees and winners from 1934 through 2014.

### Usage

```
data(best_picture)
```

### Format

A `data.frame` object with 484 observations and 19 columns. The columns are defined as follows:

`year` The integer year of the nomination. These span from 1934 through 2014. Note that the number of films nominated per year varies from 5 to 12.

`film` The title of the film.

`winner` A logical for whether the film won the Oscar for Best Picture. There is exactly one winning film per year.

`nominated_for_Writing` A logical indicating whether the film was also nominated for a Writing award that year.

`nominated_for_BestDirector` A logical indicating whether the film was also nominated for Best Director award that year.

`nominated_for_BestActress` A logical indicating whether the film was also nominated for at least one Best Actress award that year.

`nominated_for_BestActor` A logical indicating whether the film was also nominated for at least one Best Actor award that year.

`nominated_for_BestFilmEditing` A logical indicating whether the film was also nominated for at least one Best Film Editing award that year.

`Adventure` A double computed as a 0/1 indicator of whether “Adventure” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

`Biography` A double computed as a 0/1 indicator of whether “Biography” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

`Comedy` A double computed as a 0/1 indicator of whether “Comedy” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

`Crime` A double computed as a 0/1 indicator of whether “Crime” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

`Drama` A double computed as a 0/1 indicator of whether “Drama” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

`History` A double computed as a 0/1 indicator of whether “History” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

**Musical** A double computed as a 0/1 indicator of whether “Musical” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

**Romance** A double computed as a 0/1 indicator of whether “Romance” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

**Thriller** A double computed as a 0/1 indicator of whether “Thriller” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

**War** A double computed as a 0/1 indicator of whether “War” was one of the genres tagged for the film in IMDb, divided by the total count of genres tagged for the film.

**Other** A double computed as 1 minus the sum of the other genre indicators. Effectively this is the sum of indicators for “Mystery”, “Family”, “Fantasy”, “Action”, “Western”, “Music”, “Sport”, “Sci Fi”, “Film-Noir”, “Animation”, and “Horror” divided by the total count of genres tagged for the film.

### Note

“Oscar” is a copyright property of the Academy of Motion Picture Arts and Sciences. IMDb is owned by Amazon.

### Author(s)

Steven E. Pav <shabbychef@gmail.com>

### Source

Awards data were sourced from Wikipedia, while genre data were sourced from IMDb. Any errors in transcription are the fault of the package author.

### Examples

```
library(dplyr)
data(best_picture)

best_picture %>%
  group_by(nominated_for_BestDirector) %>%
  summarize(propwin=mean(winner)) %>%
  ungroup()
best_picture %>%
  group_by(nominated_for_BestActor) %>%
  summarize(propwin=mean(winner)) %>%
  ungroup()
# hmmm.
best_picture %>%
  group_by(nominated_for_BestActress) %>%
  summarize(propwin=mean(winner)) %>%
  ungroup()
```

---

diving

*Olympic Diving Data*

---

**Description**

One hundred years of Men's Olympic Platform Diving records.

**Usage**

```
data(diving)
```

**Format**

A data.frame object with 695 observations and 13 columns.

The columns are defined as follows:

Name The participant's name.

Age The age of the participant at the time of the Olympics. Some values missing.

Height The height of the participant at the time of the Olympics, in centimeters. Many values missing.

Weight The height of the participant at the time of the Olympics, in kilograms. Many values missing.

Team The string name of the team (country) which the participant represented.

NOC The string name of the National Olympic Committee which the participant represented. This is a three character code.

Games The string name of the Olympic games, including a year.

Year The integer year of the Olympics. These range from 1906 through 2016.

City The string name of the host city.

Medal A string of "Gold", "Silver", "Bronze" or NA.

EventId A unique integer ID for each Olympics.

AthleteId A unique integer ID for each participant.

HOST\_NOC The string name of the National Olympic Committee of the nation hosting the Olympics. This is a three character code.

**Note**

The author makes no guarantees regarding correctness of this data.

Please attribute this data to the upstream harvester.

**Author(s)**

Steven E. Pav <shabbychef@gmail.com>

**Source**

Data were collected by Randi Griffin from the website “sports-reference.com”, and staged on Kaggle at <https://www.kaggle.com/heesoo37/120-years-of-olympic-history-athletes-and-results>.

**Examples**

```
library(dplyr)
library(forcats)
data(diving)

fitdat <- diving %>%
  mutate(Finish=case_when(grepl('Gold',Medal) ~ 1,
                           grepl('Silver',Medal) ~ 2,
                           grepl('Bronze',Medal) ~ 3,
                           TRUE ~ 4)) %>%
  mutate(weight=ifelse(Finish <= 3,1,0)) %>%
  mutate(cut_age=cut(coalesce(Age,22.0),c(12,19.5,21.5,22.5,25.5,99),include.lowest=TRUE)) %>%
  mutate(country=forcats::fct_relevel(forcats::fct_lump(factor(NOC),n=5),'Other')) %>%
  mutate(home_advantage=NOC==HOST_NOC)

hensm(Finish ~ cut_age + country + home_advantage,data=fitdat,weights=weight,group=EventId,ngamma=3)
```

erank

*Expected rank under the Harville model.***Description**

Compute the expected rank of a bunch of entries based on their probability of winning under the Harville model.

**Usage**

```
erank(mu)
```

**Arguments**

mu a vector giving the probabilities. Should sum to one.

**Details**

Given the vector  $\mu$ , we compute the expected rank of each entry, under the Harville model, where tail probabilities of winning remain proportional.

Under the Harville model, the probability that the  $i$ th element is assigned value 1 is

$$\pi_{1,i} = \frac{\mu_i}{\sum_j \mu_j}.$$

Once an element has been assigned a 1, the Harville procedure removes it from the set and iterates. If there are  $k$  elements in  $\mu$ , then the  $i$ th element can be assigned any place between 1 and  $k$ . This function computes the expected value of that random variable.

While a naive implementation of this function would take time factorial in  $k$ , this implementation takes time quadratic in  $k$ , since it can be shown that the expected rank of the  $i$ th element takes value

$$e_i = k + \frac{1}{2} - \sum_j \frac{\mu_i}{\mu_i + \mu_j}.$$

### Value

The expected ranks, a vector.

### Note

we should have the sum of ranks equal to the sum of  $1 : \text{length}(\mu)$ .

### Author(s)

Steven E. Pav <shabbychef@gmail.com>

### Examples

```
# a garbage example
set.seed(12345)
mus <- runif(12)
mus <- mus / sum(mus)
erank(mus)

# confirm the expected rank via simulation
set.seed(123)
mus <- runif(6,min=0,max=2)
mus <- mus / sum(mus)
set.seed(101)
emp <- rowMeans(replicate(200,rhenery(mu=mus,gamma=rep(1,length(mus)-1))))
(emp - erank(mus)) / emp

if (require(microbenchmark)) {
  p10 <- 1:10 / sum(1:10)
  p16 <- 1:16 / sum(1:16)
  p24 <- 1:24 / sum(1:24)
  microbenchmark(erank(p10), erank(p16), erank(p24))
}
```

---

`glance.harsm`*Glance*

---

**Description**

Glance at a harsm, hensm or linodds model.

**Usage**

```
## S3 method for class 'harsm'  
glance(x, ...)  
  
## S3 method for class 'hensm'  
glance(x, ...)
```

**Arguments**

`x` an object of type harsm or hensm  
`...` arguments for generic consistency.

**Details**

Returns a table with information on the overall fit an estimated Harville or Henery model.

**Value**

A glanced `tibble::tibble()` with fields

**df** The degrees of freedom of the model.

**logLik** The log-likelihood of the model.

**AIC** Akaike's Information Criterion for the model.

**nobs** The number of observations, if this is available, otherwise 'NA'.

**Note**

In the future this may include information about the regularization, if any.

**Author(s)**

Steven E. Pav <shabbychef@gmail.com>

**See Also**

[glance.maxLik](#).

**Examples**

```
# softmax on the Best Picture data
data(best_picture)
df <- best_picture
df$place <- ifelse(df$winner,1,2)
df$weight <- ifelse(df$winner,1,0)

fmla <- place ~ nominated_for_BestDirector + nominated_for_BestActor + Drama
fit0 <- harsm(fmla,data=df,group=year,weights=weight)
print(glance(fit0))
```

---

harsm

*Friendly interface to softmax regression under Harville model.*


---

**Description**

A user friendly interface to the softmax regression under the Harville model.

**Usage**

```
harsm(
  formula,
  data,
  group = NULL,
  weights = NULL,
  fit0 = NULL,
  reg_wt = NULL,
  reg_zero = 0,
  reg_power = NULL,
  reg_coef_idx = NULL,
  reg_standardize = FALSE,
  na.action = na.omit
)

## S3 method for class 'harsm'
vcov(object, ...)

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

**Arguments**

formula            an object of class "[formula](#)" (or one that can be coerced to that class): a symbolic description of the model to be fitted. The details of model specification are given under 'Details'.

data	an optional data frame, list or environment (or object coercible by <code>as.data.frame</code> to a data frame) containing the variables in the model. If not found in data, the variables are taken from <code>environment(formula)</code> , typically the environment from which <code>lm</code> is called.
group	the string name of the group variable in the data, or a bare character with the group name. The group indices need not be integers, but that is more efficient. They need not be sorted.
weights	an optional vector of weights, or the string or bare name of the weights in the data for use in the fitting process. The weights are attached to the outcomes, not the participant. Set to NULL for none.
fit0	An optional object of class <code>harsm</code> or of <code>hensm</code> with the initial fit estimates. These will be used for ‘warm start’ of the estimation procedure. A warm start should only speed up estimation, not change the ultimate results. When there is mismatch between the coefficients in <code>fit0</code> and the model being fit here, the missing coefficients are initialized as zero.
reg_wt	the multiplicative weight(s) of the regularization terms. must be non-negative. May be a scalar or vector, but will be recycled to the length of the <code>reg_coef_idx</code> .
reg_zero	the ‘zero’ of the regularization terms. This would usually be zero if you want to shrink to zero, but in some cases you may wish to shrink to 1 for example. May be a scalar or vector, but will be recycled to the length of the <code>reg_coef_idx</code> . If NULL is given, defaults to zeroes for beta terms, and ones for gamma terms in a Henery model fit, all zeroes for Harville model fits.
reg_power	the power of the regularization terms, 2 for ridge regression, 1 for lasso.
reg_coef_idx	the index of the coefficient which the corresponding regularization term is applied to. For the Harville model, the indices only refer to the $\beta$ coefficient vector. For the Henery model, the indices refer to the <i>beta</i> coefficient vector and $\gamma$ coefficient vector concatenated together. The other regularization parameters may be recycled as needed, but not the <code>reg_coef_idx</code> values.
reg_standardize	if true, the <code>reg_wt</code> are normalized, or ‘standardized’ with respect to the standard deviation of the corresponding columns of the design matrix. That is, the weight used is the given weight <i>times</i> the standard deviation of the corresponding independent variable to the corresponding <code>reg_power</code> . Only terms associated with the betas are so normalized.
na.action	How to deal with missing values in the outcomes, groups, weights, etc.
object	an object of class <code>harsm</code> .
...	For <code>lm()</code> : additional arguments to be passed to the low level regression fitting functions (see below).
x	an object used to select a method.

## Details

Performs a softmax regression by groups, via Maximum Likelihood Estimation, under the Harville model. We fit  $\beta$  where odds are  $\eta = x^T \beta$  for independent variables  $x$ . The probability of taking first place is then  $\mu = c \exp \eta$ , where the  $c$  is chosen so the  $\mu$  sum to one. Under the Harville model,

conditional on the first place finisher having been observed, the probability model for second (and successive) places with the probabilities of the remaining participants renormalized.

The `print` method of the `harsm` object includes a display of the R-squared. This measures the improvement in squared errors of the expected rank from the model over the null model which posits that all odds are equal. When the formula includes an offset, a ‘delta R-squared’ is also output. This is the improvement in predicted ranks over the model based on the offset term. Note that the expected ranks are only easy to produce under the Harville model; under the Henery model, the summary R-squared statistics are not produced. Note that this computation uses weighted sums of squares, as the weights contribute to the likelihood term. However, the square sum computation does not take into account the standard error of the rank, and so unlike in linear regression, the softmax regression does not always give positive R-squareds, and the statistic is otherwise hard to interpret.

### Value

An object of class `harsm`, but also of `maxLik` with the fit.

### Regularization

The regularization term is of the form

$$\sum_i w_i |\nu_{c_i} - z_i|^{p_i},$$

where  $w_i$  are the `reg_wt` weights,  $z_i$  are the `reg_zero` zeroes,  $p_i$  are the `reg_power` powers, and  $c_i$  are the `reg_coef_idx` coefficient indices. Note that the coefficient indices can be repeated so that the regularization term can include multiple contributions from one coefficient element. This allows ‘elasticnet’ regularization. The  $\nu$  here refer to the regression coefficients  $\beta$  concatenated with the  $\gamma$  coefficients in the Henery model.

### Note

The fit functions return an object of type `maxLik` even when regularization penalties are applied; the statistical inference functions are not valid when regularization is used. The user is warned.

This regression may give odd results when the outcomes are tied, imposing an arbitrary order on the tied outcomes. Moreover, no warning may be issued in this case. In future releases, ties may be dealt with differently, perhaps in analogy to how ties are treated in the Cox Proportional Hazards regression, using the methods of Breslow or Efron.

To avoid incorrect inference when only the top performers are recorded, and all others are effectively tied, one should use weighting. Set the weights to zero for participants who are tied non-winners, and one for the rest. So for example, if you observe the Gold, Silver, and Bronze medal winners of an Olympic event that had a starting field of 12 participants, set weights to 1 for the medal winners, and 0 for the others. Note that the weights do not attach to the participants, they attach to the place they took.

Since version 0.1.0 of this package, the normalization of weights used in this function have changed under the hood. This is to give correct inference in the case where zero weights are used to signify finishing places were not observed. If in doubt, please confirm inference by simulations, taking as example the simulations in the README.

**Author(s)**

Steven E. Pav <shabbychef@gmail.com>

**See Also**

[harsmfit](#), [harsmlik](#).

**Examples**

```
nfeat <- 5
set.seed(1234)
g <- ceiling(seq(0.1,1000,by=0.1))
X <- matrix(rnorm(length(g) * nfeat),ncol=nfeat)
beta <- rnorm(nfeat)
eta <- X %%% beta
y <- rsm(eta,g)
# now the pretty frontend
data <- cbind(data.frame(outcome=y,race=g),as.data.frame(X))

fmla <- outcome ~ V1 + V2 + V3 + V4 + V5
fitm <- harsm(fmla,data,group=race)

eta0 <- rowMeans(X)
data <- cbind(data.frame(outcome=y,race=g,eta0=eta0),as.data.frame(X))
fmla <- outcome ~ offset(eta0) + V1 + V2 + V3 + V4 + V5
fitm <- harsm(fmla,data,group=race)

# with weights
data <- cbind(data.frame(outcome=y,race=g,eta0=eta0),as.data.frame(X))
data$wts <- runif(nrow(data),min=1,max=2)
fmla <- outcome ~ offset(eta0) + V1 + V2 + V3 + V4 + V5
fitm <- harsm(fmla,data,group=race,weights=wts)

# softmax on the Best Picture data
data(best_picture)
df <- best_picture
df$place <- ifelse(df$winner,1,2)
df$weight <- ifelse(df$winner,1,0)

fmla <- place ~ nominated_for_BestDirector + nominated_for_BestActor + Drama
fit0 <- harsm(fmla,data=df,group=year,weights=weight)

# warm start is a thing:
sub_fmla <- place ~ nominated_for_BestDirector + nominated_for_BestActor
fit1 <- harsm(sub_fmla,data=df,group=year,weights=weight,fit0=fit0)

# ridge regression.
fitr2 <- harsm(sub_fmla,data=df,group=year,weights=weight,
reg_wt=rep(1,2), reg_power=2, reg_zero=0, reg_coef_idx=c(1,2))

# l1 regularization regression.
fitr1 <- harsm(sub_fmla,data=df,group=year,weights=weight,
```

```

reg_wt=rep(1,2), reg_power=1, reg_zero=0, reg_coef_idx=c(1,2))

# elasticnet regularization regression.
fitr12 <- harsm(sub_fm1a,data=df,group=year,weights=weight,
reg_wt=rep(1,4), reg_power=c(1,1,2,2), reg_zero=0, reg_coef_idx=c(1,2,1,2))

# test against logistic regression
if (require(dplyr)) {
nevent <- 10000
set.seed(1234)
adf <- data_frame(eventnum=floor(seq(1,nevent + 0.7,by=0.5))) %>%
  mutate(x=rnorm(n()),
         program_num=rep(c(1,2),nevent),
         intercept=as.numeric(program_num==1),
         eta=1.5 * x + 0.3 * intercept,
         place=ohenerly::rsm(eta,g=eventnum))

# Harville model
modh <- harsm(place ~ intercept + x,data=adf,group=eventnum)

# the collapsed data.frame for glm
ddf <- adf %>%
  arrange(eventnum,program_num) %>%
  group_by(eventnum) %>%
  summarize(resu=as.numeric(first(place)==1),
           delx=first(x) - last(x),
           deli=first(intercept) - last(intercept)) %>%
  ungroup()

# glm logistic fit
modg <- glm(resu ~ delx + 1,data=ddf,family=binomial(link='logit'))

all.equal(as.numeric(coef(modh)),as.numeric(coef(modg)),tolerance=1e-4)
all.equal(as.numeric(vcov(modh)),as.numeric(vcov(modg)),tolerance=1e-4)

}

```

---

harsmfit

*Experts only softmax regression under Harville model.*


---

### Description

An “experts only” softmax fitting function for the Harville model.

**Usage**

```

harsmfit(
  y,
  g,
  X,
  wt = NULL,
  eta0 = NULL,
  beta0 = NULL,
  normalize_wt = FALSE,
  reg_wt = NULL,
  reg_zero = 0,
  reg_power = NULL,
  reg_coef_idx = NULL,
  reg_standardize = FALSE,
  method = c("BFGS", "NR", "CG", "NM")
)

```

**Arguments**

<code>y</code>	a vector of the ranked outcomes within each group. Only the order within a group matters.
<code>g</code>	a vector giving the group indices. Need not be integers, but that is more efficient. Need not be sorted. Must be the same length as <code>y</code> .
<code>X</code>	a matrix of the independent variables. Must have as many rows as the length of <code>y</code> .
<code>wt</code>	an optional vector of the observation level weights. These must be non-negative, otherwise an error is thrown. Note that the weight of the last ranked outcome within a group is essentially ignored. Must be the same length as <code>y</code> .
<code>eta0</code>	an optional vector of the consensus odds. These are added to the fit odds in odds space before the likelihood calculation. If given, then when the model is used to predict, similar consensus odds must be given. Must be the same length as <code>y</code> .
<code>beta0</code>	an optional vector of the initial estimate of beta for ‘warm start’ of the estimation procedure. Must be the same length as number of columns in <code>X</code> . Should only affect the speed of the computation, not the results. Defaults to all zeroes.
<code>normalize_wt</code>	if TRUE, we renormalize <code>wt</code> , if given, to have mean value 1. Note that the default value has changed since version 0.1.0 of this package. Moreover, non-normalized weights can lead to incorrect inference. Use with caution.
<code>reg_wt</code>	the multiplicative weight(s) of the regularization terms. must be non-negative. May be a scalar or vector, but will be recycled to the length of the <code>reg_coef_idx</code> .
<code>reg_zero</code>	the ‘zero’ of the regularization terms. This would usually be zero if you want to shrink to zero, but in some cases you may wish to shrink to 1 for example. May be a scalar or vector, but will be recycled to the length of the <code>reg_coef_idx</code> . If NULL is given, defaults to zeroes for beta terms, and ones for gamma terms in a Henry model fit, all zeroes for Harville model fits.
<code>reg_power</code>	the power of the regularization terms, 2 for ridge regression, 1 for lasso.

reg_coef_idx	the index of the coefficient which the corresponding regularization term is applied to. For the Harville model, the indices only refer to the $\beta$ coefficient vector. For the Henery model, the indices refer to the <i>beta</i> coefficient vector and $\gamma$ coefficient vector concatenated together. The other regularization parameters may be recycled as needed, but not the reg_coef_idx values.
reg_standardize	if true, the reg_wt are normalized, or ‘standardized’ with respect to the standard deviation of the corresponding columns of the design matrix. That is, the weight used is the given weight <i>times</i> the standard deviation of the corresponding independent variable to the corresponding reg_power. Only terms associated with the betas are so normalized.
method	maximisation method, currently either "NR" (for Newton-Raphson), "BFGS" (for Broyden-Fletcher-Goldfarb-Shanno), "BFGSR" (for the BFGS algorithm implemented in R), "BHHH" (for Berndt-Hall-Hall-Hausman), "SANN" (for Simulated ANNealing), "CG" (for Conjugate Gradients), or "NM" (for Nelder-Mead). Lower-case letters (such as "nr" for Newton-Raphson) are allowed. The default method is "NR" for unconstrained problems, and "NM" or "BFGS" for constrained problems, depending on if the grad argument was provided. "BHHH" is a good alternative given the likelihood is returned observation-wise (see <a href="#">maxBHHH</a> ). Note that stochastic gradient ascent (SGA) is currently not supported as this method seems to be rarely used for maximum likelihood estimation.

### Details

Given a number of events, indexed by group, and a vector  $y$  of the ranks of each entry within that group, perform maximum likelihood estimation under the softmax and proportional probability model.

The user can optionally supply a vector of  $\eta_0$ , which are taken as the fixed, or ‘consensus’ odds. The estimation is then conditional on these fixed odds.

Weighted estimation is supported.

The code relies on the likelihood function of [harsmlik](#), and MLE code from [maxLik](#).

### Value

An object of class `harsm`, `maxLik`, and `linodds`.

### Regularization

The regularization term is of the form

$$\sum_i w_i |\nu_{c_i} - z_i|^{p_i},$$

where  $w_i$  are the reg\_wt weights,  $z_i$  are the reg\_zero zeroes,  $p_i$  are the reg\_power powers, and  $c_i$  are the reg\_coef\_idx coefficient indices. Note that the coefficient indices can be repeated so that the regularization term can include multiple contributions from one coefficient element. This allows ‘elasticnet’ regularization. The  $\nu$  here refer to the regression coefficients  $\beta$  concatenated with the  $\gamma$  coefficients in the Henery model.

**Note**

The fit functions return an object of type `maxLik` even when regularization penalties are applied; the statistical inference functions are not valid when regularization is used. The user is warned.

**Author(s)**

Steven E. Pav <shabbychef@gmail.com>

**References**

Harville, D. A. "Assigning probabilities to the outcomes of multi-entry competitions." *Journal of the American Statistical Association* 68, no. 342 (1973): 312-316. doi:10.1080/01621459.1973.10482425

**See Also**

the friendly interface `harsm`, the likelihood function, `harsmlik`, and the expected rank function (the inverse link), `erank`.

**Examples**

```
nfeat <- 5
set.seed(1234)
g <- ceiling(seq(0.1,1000,by=0.1))
X <- matrix(rnorm(length(g) * nfeat),ncol=nfeat)
beta <- rnorm(nfeat)
eta <- X %*% beta
y <- rsm(eta,g)

mod0 <- harsmfit(y=y,g=g,X=X)
summary(mod0)
# now upweight finishers 1-5
modw <- harsmfit(y=y,g=g,X=X,wt=1 + as.numeric(y < 6))
summary(modw)
```

---

harsm\_invlink

*The inverse link for the softmax.*

---

**Description**

The inverse link function for the softmax. This function takes the group-wise probabilities,  $\mu$ , and computes the expected ranks within each group under the Harville model. That is, it is a groupwise computation of the `erank` function.

**Usage**

```
harsm_invlink(eta, mu = smax(eta, g), g = NULL)
```

**Arguments**

eta	a vector of the odds. Must be the same length as g if g is given.
mu	a vector of the probabilities. Should sum to one, at least per group. Should be the same size as g if given. If both mu and eta are given, a warning is issued, and the mu is used.
g	a vector giving the group indices. If NULL, then we assume only one group is in consideration.

**Value**

a vector of the ranks.

**Author(s)**

Steven E. Pav <shabbychef@gmail.com>

**See Also**

the ungrouped version of this, [erank](#).

**Examples**

```
mus <- runif(12)
mus <- mus / sum(mus)
harsm_invlink(mus)

harsm_invlink(mus, c(rep(1,6), rep(2,6)))
```

---

hensm

*Friendly interface to softmax regression under Henery model.*

---

**Description**

A user friendly interface to the softmax regression under the Henery model.

**Usage**

```
hensm(
  formula,
  data,
  group = NULL,
  weights = NULL,
  ngamma = 4,
  fit0 = NULL,
  reg_wt = NULL,
  reg_zero = NULL,
  reg_power = NULL,
```

```

    reg_coef_idx = NULL,
    reg_standardize = FALSE,
    na.action = na.omit
  )

## S3 method for class 'hensm'
vcov(object, ...)

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

```

### Arguments

formula	an object of class " <a href="#">formula</a> " (or one that can be coerced to that class): a symbolic description of the model to be fitted. The details of model specification are given under 'Details'.
data	an optional data frame, list or environment (or object coercible by <a href="#">as.data.frame</a> to a data frame) containing the variables in the model. If not found in data, the variables are taken from <code>environment(formula)</code> , typically the environment from which <code>lm</code> is called.
group	the string name of the group variable in the data, or a bare character with the group name. The group indices need not be integers, but that is more efficient. They need not be sorted.
weights	an optional vector of weights, or the string or bare name of the weights in the data for use in the fitting process. The weights are attached to the outcomes, not the participant. Set to <code>NULL</code> for none.
ngamma	The number of gammas to fit. Should be at least 2.
fit0	An optional object of class <code>hensm</code> or of <code>harsm</code> with the initial fit estimates. These will be used for 'warm start' of the estimation procedure. A warm start should only speed up estimation, not change the ultimate results. When there is mismatch between the coefficients in <code>fit0</code> and the model being fit here, the missing coefficients are initialized as zero. If <code>ngamma</code> is <code>NULL</code> and <code>fit0</code> is given, we default to the number of gammas in the initial fit, otherwise we fill any missing gammas with 1. If a <code>harsm</code> object is given, then <code>ngamma</code> must be non-null.
reg_wt	the multiplicative weight(s) of the regularization terms. must be non-negative. May be a scalar or vector, but will be recycled to the length of the <code>reg_coef_idx</code> .
reg_zero	the 'zero' of the regularization terms. This would usually be zero if you want to shrink to zero, but in some cases you may wish to shrink to 1 for example. May be a scalar or vector, but will be recycled to the length of the <code>reg_coef_idx</code> . If <code>NULL</code> is given, defaults to zeroes for beta terms, and ones for gamma terms in a Henery model fit, all zeroes for Harville model fits.
reg_power	the power of the regularization terms, 2 for ridge regression, 1 for lasso.
reg_coef_idx	the index of the coefficient which the corresponding regularization term is applied to. For the Harville model, the indices only refer to the $\beta$ coefficient vector. For the Henery model, the indices refer to the <i>beta</i> coefficient vector and $\gamma$ coefficient vector concatenated together. The other regularization parameters may be recycled as needed, but not the <code>reg_coef_idx</code> values.

reg_standardize	if true, the reg_wt are normalized, or ‘standardized’ with respect to the standard deviation of the corresponding columns of the design matrix. That is, the weight used is the given weight <i>times</i> the standard deviation of the corresponding independent variable to the corresponding reg_power. Only terms associated with the betas are so normalized.
na.action	How to deal with missing values in y, g, X, wt, eta0.
object	an object of class hensm.
...	For lm(): additional arguments to be passed to the low level regression fitting functions (see below).
x	an object used to select a method.

### Details

Performs a softmax regression by groups, via Maximum Likelihood Estimation. It is assumed that successive sub-races maintain the proportional probability of the softmax, up to some gamma coefficients,  $\gamma_2, \gamma_3, \dots, \gamma_n$ , which we fit. This model nests the Harville model fit by `harsm`, by fixing all the gammas equal to 1.

### Value

An object of class `hensm`, but also of `maxLik` with the fit.

### Regularization

The regularization term is of the form

$$\sum_i w_i |\nu_{c_i} - z_i|^{p_i},$$

where  $w_i$  are the reg\_wt weights,  $z_i$  are the reg\_zero zeroes,  $p_i$  are the reg\_power powers, and  $c_i$  are the reg\_coef\_idx coefficient indices. Note that the coefficient indices can be repeated so that the regularization term can include multiple contributions from one coefficient element. This allows ‘elasticnet’ regularization. The  $\nu$  here refer to the regression coefficients  $\beta$  concatenated with the  $\gamma$  coefficients in the Henery model.

### Note

The fit functions return an object of type `maxLik` even when regularization penalties are applied; the statistical inference functions are not valid when regularization is used. The user is warned.

This regression may give odd results when the outcomes are tied, imposing an arbitrary order on the tied outcomes. Moreover, no warning may be issued in this case. In future releases, ties may be dealt with differently, perhaps in analogy to how ties are treated in the Cox Proportional Hazards regression, using the methods of Breslow or Efron.

To avoid incorrect inference when only the top performers are recorded, and all others are effectively tied, one should use weighting. Set the weights to zero for participants who are tied non-winners, and one for the rest. So for example, if you observe the Gold, Silver, and Bronze medal winners of an Olympic event that had a starting field of 12 participants, set weights to 1 for the medal winners,

and 0 for the others. Note that the weights do not attach to the participants, they attach to the place they took.

Since version 0.1.0 of this package, the normalization of weights used in this function have changed under the hood. This is to give correct inference in the case where zero weights are used to signify finishing places were not observed. If in doubt, please confirm inference by simulations, taking as example the simulations in the README.

When regularization is used, the first gamma coefficient is not shrunk, as it always equals one in the Henery model, and is not estimated from the data.

### Author(s)

Steven E. Pav <shabbychef@gmail.com>

### See Also

[harsm](#), [smlk](#).

### Examples

```
nfeat <- 5
set.seed(1234)
g <- ceiling(seq(0.1,1000,by=0.1))
X <- matrix(rnorm(length(g) * nfeat),ncol=nfeat)
beta <- rnorm(nfeat)
eta <- X %*% beta
# 2FIX: do rhenery
y <- rsm(eta,g)
# now the pretty frontend
data <- cbind(data.frame(outcome=y,race=g),as.data.frame(X))

fmla <- outcome ~ V1 + V2 + V3 + V4 + V5
fitm <- hensm(fmla,data,group=race)

# with offset
eta0 <- rowMeans(X)
data <- cbind(data.frame(outcome=y,race=g,eta0=eta0),as.data.frame(X))
fmla <- outcome ~ offset(eta0) + V1 + V2 + V3 + V4 + V5
fitm <- hensm(fmla,data,group=race)

# on horse race data
library(dplyr)
data(race_data)
df <- race_data %>%
  group_by(EventId) %>%
  mutate(eta0=log(WN_pool / sum(WN_pool))) %>%
  ungroup() %>%
  mutate(weights=ifelse(!is.na(Finish),1,0)) %>%
  mutate(fac_age=cut(Age,c(0,3,5,7,Inf),include.lowest=TRUE))

# Henery Model with market efficiency
hensm(Finish ~ eta0,data=df,group=EventId,weights=weights,ngamma=3)
```

```

# look for age effect not captured by consensus odds.
fmla <- Finish ~ offset(eta0) + fac_age
fit0 <- hensm(fmla,data=df,group=EventId,weights=weights,ngamma=2)
# allow warm start.
fit1 <- hensm(fmla,data=df,group=EventId,weights=weights,fit0=fit0,ngamma=2)
# allow warm start with more gammas.
fit2 <- hensm(fmla,data=df,group=EventId,weights=weights,fit0=fit0,ngamma=3)
# or a different formula
fit3 <- hensm(update(fmla,~ . + PostPosition),data=df,group=EventId,weights=weights,fit0=fit0)

# warm start from harsm object
fit0_har <- harsm(fmla,data=df,group=EventId,weights=weights)
fit4 <- hensm(fmla,data=df,group=EventId,fit0=fit0_har,weights=weights)

# regularization examples
fmla <- Finish ~ offset(eta0) + fac_age + PostPosition
# no regularization
fitr0 <- hensm(fmla,data=df,group=EventId,weights=weights,ngamma=3)
# ridge regression on the betas (there are two)
fitr2 <- hensm(fmla,data=df,group=EventId,weights=weights,ngamma=3,
  reg_wt=rep(1,2), reg_power=2, reg_zero=0, reg_coef_idx=c(1,2))
# l1 regression on the betas (there are two)
fitr1 <- hensm(fmla,data=df,group=EventId,weights=weights,ngamma=3,
  reg_wt=rep(1,2), reg_power=1, reg_zero=0, reg_coef_idx=c(1,2))
# elasticnet regression on the betas (there are two)
fitr12 <- hensm(fmla,data=df,group=EventId,weights=weights,ngamma=3,
  reg_wt=rep(1,4), reg_power=c(1,1,2,2), reg_zero=0, reg_coef_idx=c(1,2,1,2))
# l1 regression on the gammas, shrinking to one (there are two estimated)
fitrg1 <- hensm(fmla,data=df,group=EventId,weights=weights,ngamma=3,
  reg_wt=rep(1,2), reg_power=1, reg_zero=1, reg_coef_idx=c(3,4))
# l2 regression on the betas and gammas, shrinking beta to 0, gamma to 1.
fitrg2 <- hensm(fmla,data=df,group=EventId,weights=weights,ngamma=3,
  reg_wt=rep(1,4), reg_power=2, reg_zero=c(0,0,1,1), reg_coef_idx=1:4)

```

---

hensmfit

*Experts only softmax regression under Henery model.*


---

## Description

An “experts only” softmax fitting function for the Henery model.

## Usage

```

hensmfit(
  y,
  g,
  X,
  wt = NULL,

```

```

eta0 = NULL,
beta0 = NULL,
gamma0 = NULL,
normalize_wt = FALSE,
ngamma = 4,
reg_wt = NULL,
reg_zero = NULL,
reg_power = NULL,
reg_coef_idx = NULL,
reg_standardize = FALSE,
method = c("BFGS", "NR", "CG", "NM")
)

```

### Arguments

y	a vector of the ranked outcomes within each group. Only the order within a group matters.
g	a vector giving the group indices. Need not be integers, but that is more efficient. Need not be sorted. Must be the same length as y.
X	a matrix of the independent variables. Must have as many rows as the length of y.
wt	an optional vector of the observation level weights. These must be non-negative, otherwise an error is thrown. Note that the weight of the last ranked outcome within a group is essentially ignored. Must be the same length as y.
eta0	an optional vector of the consensus odds. These are added to the fit odds in odds space before the likelihood calculation. If given, then when the model is used to predict, similar consensus odds must be given. Must be the same length as y.
beta0	an optional vector of the initial estimate of beta for 'warm start' of the estimation procedure. Must be the same length as number of columns in X. Should only affect the speed of the computation, not the results. Defaults to all zeroes.
gamma0	an optional vector of the initial estimate of beta for 'warm start' of the estimation procedure. Must be of length $ngamma - 1$ . Should only affect the speed of the computation, not the results. Defaults to all ones.
normalize_wt	if TRUE, we renormalize wt, if given, to have mean value 1. Note that the default value has changed since version 0.1.0 of this package. Moreover, non-normalized weights can lead to incorrect inference. Use with caution.
ngamma	The number of gammas to fit. Should be at least 2.
reg_wt	the multiplicative weight(s) of the regularization terms. must be non-negative. May be a scalar or vector, but will be recycled to the length of the reg_coef_idx.
reg_zero	the 'zero' of the regularization terms. This would usually be zero if you want to shrink to zero, but in some cases you may wish to shrink to 1 for example. May be a scalar or vector, but will be recycled to the length of the reg_coef_idx. If NULL is given, defaults to zeroes for beta terms, and ones for gamma terms in a Henry model fit, all zeroes for Harville model fits.
reg_power	the power of the regularization terms, 2 for ridge regression, 1 for lasso.

reg_coef_idx	the index of the coefficient which the corresponding regularization term is applied to. For the Harville model, the indices only refer to the $\beta$ coefficient vector. For the Henery model, the indices refer to the <i>beta</i> coefficient vector and $\gamma$ coefficient vector concatenated together. The other regularization parameters may be recycled as needed, but not the reg_coef_idx values.
reg_standardize	if true, the reg_wt are normalized, or ‘standardized’ with respect to the standard deviation of the corresponding columns of the design matrix. That is, the weight used is the given weight <i>times</i> the standard deviation of the corresponding independent variable to the corresponding reg_power. Only terms associated with the betas are so normalized.
method	maximisation method, currently either "NR" (for Newton-Raphson), "BFGS" (for Broyden-Fletcher-Goldfarb-Shanno), "BFGSR" (for the BFGS algorithm implemented in R), "BHHH" (for Berndt-Hall-Hall-Hausman), "SANN" (for Simulated ANNealing), "CG" (for Conjugate Gradients), or "NM" (for Nelder-Mead). Lower-case letters (such as "nr" for Newton-Raphson) are allowed. The default method is "NR" for unconstrained problems, and "NM" or "BFGS" for constrained problems, depending on if the grad argument was provided. "BHHH" is a good alternative given the likelihood is returned observation-wise (see <a href="#">maxBHHH</a> ).  Note that stochastic gradient ascent (SGA) is currently not supported as this method seems to be rarely used for maximum likelihood estimation.

## Details

Given a number of events, indexed by group, and a vector  $y$  of the ranks of each entry within that group, perform maximum likelihood estimation under the softmax and proportional probability model.

The user can optionally supply a vector of  $\eta_0$ , which are taken as the fixed, or ‘consensus’ odds. The estimation is then conditional on these fixed odds.

Weighted estimation is supported.

The code relies on the likelihood function of [hensmlik](#), and MLE code from [maxLik](#).

## Value

An object of class `hensm`, `maxLik`, and `linodds`.

## Author(s)

Steven E. Pav <[shabbychef@gmail.com](mailto:shabbychef@gmail.com)>

## See Also

the friendly interface `hensm`, the likelihood function, [hensmlik](#), and the expected rank function (the inverse link), [erank](#).

**Examples**

```

nfeat <- 5
set.seed(1234)
g <- ceiling(seq(0.1,1000,by=0.1))
X <- matrix(rnorm(length(g) * nfeat),ncol=nfeat)
beta <- rnorm(nfeat)
eta <- X %*% beta
y <- rsm(eta,g)

mod0 <- hensmfit(y=y,g=g,X=X,ngamma=3)
summary(mod0)
# now upweight finishers 1-5
modw <- hensmfit(y=y,g=g,X=X,wt=1 + as.numeric(y < 6),ngamma=2)
summary(modw)

```

---

inv\_smax

*The inverse softmax function.*


---

**Description**

The inverse softmax function: take a logarithm and center.

**Usage**

```
inv_smax(mu, g = NULL)
```

**Arguments**

mu	a vector of the probabilities. Must be the same length as g if g is given. If mu and eta are both given, we ignore eta and use mu.
g	a vector giving the group indices. If NULL, then we assume only one group is in consideration.

**Details**

This is the inverse of the softmax function. Given vector  $\mu$  for a single group, finds vector  $\eta$  such that

$$\eta_i = \log \mu_i + c,$$

where  $c$  is chosen such that the  $\eta$  sum to zero:

$$c = \frac{-1}{n} \sum_i \log \mu_i.$$

**Value**

the centered log probabilities.

**Note**

This function can deal with overflow in a semi-coherent way.

**Author(s)**

Steven E. Pav <shabbychef@gmail.com>

**See Also**

[smax](#)

**Examples**

```
# we can deal with large values:
set.seed(2345)
eta <- rnorm(12, sd=1000)
mu <- smax(eta)
eta0 <- inv_smax(mu)
```

---

normalize

*Normalize a vector to sum to one.*

---

**Description**

Divide a vector by its sum, resulting in a vector with sum equal to one.

**Usage**

```
normalize(x)
```

**Arguments**

x                      vector of input data.

**Value**

the input divided by its sum. For the row-wise version, each row is divided by its sum.

**Note**

This function will return NA when any elements of the input are NA. May return Inf if the elements sum to zero.

**Author(s)**

Steven E. Pav <shabbychef@gmail.com>

ohenery

*The 'ohenery' package.***Description**

Modeling of ordinal outcomes via the softmax function under the Harville and Henery models.

**Harville and Henery models**

The Harville and Henery models describe the probability of ordered outcomes in terms of some parameters. Typically the ordered outcomes are things like place in a race, or winner among a large number of contestants. The Harville model could be described as a softmax probability for the first place finish, with a recursive model on the remaining places. The Henery model generalizes that to adjust the remaining places with another parameter.

These are best illustrated with an example. Suppose you observe a race of 20 contestants. Contestant number 11 takes first place, number 6 takes second place, and 17 takes third place, while the fourth through twentieth places are not recorded or not of interest. Under the Harville model, the probability of this outcome can be expressed as

$$\frac{\mu_{11}}{\sum_i \mu_i} \frac{\mu_6}{\sum_{i \neq 11} \mu_i} \frac{\mu_{17}}{\sum_{i \neq 11, i \neq 6} \mu_i},$$

where  $\mu_i = \exp \eta_i$ . In a softmax regression under the Harville model, one expresses the odds as  $\eta_i = x_i^\top \beta$ , where  $x_i$  are independent variables, for some  $\beta$  to be fit by the regression.

Under the Henery model, one adds gammas,  $\gamma_2, \gamma_3, \dots$  such that the probability of the outcome above is

$$\frac{\mu_{11}}{\sum_i \mu_i} \frac{\mu_6^{\gamma_2}}{\sum_{i \neq 11} \mu_i^{\gamma_2}} \frac{\mu_{17}^{\gamma_3}}{\sum_{i \neq 11, i \neq 6} \mu_i^{\gamma_3}}.$$

There is no reason to model a  $\gamma_1$  as anything but one, since it would be redundant. The Henery softmax regression estimates the  $\beta$  as well as the  $\gamma_j$ . To simplify the regression, the higher order gammas are assumed to equal the last fit value. That is, we usually model  $\gamma_5 = \gamma_4 = \gamma_3$ .

The regression supports weighted estimation as well. The weights are applied to the *places*, not to the participants. The weighted likelihood under the example above, for the Harville model is

$$\left( \frac{\mu_{11}}{\sum_i \mu_i} \right)^{w_1} \left( \frac{\mu_6}{\sum_{i \neq 11} \mu_i} \right)^{w_2} \left( \frac{\mu_{17}}{\sum_{i \neq 11, i \neq 6} \mu_i} \right)^{w_3}.$$

The weighting mechanism is how this package deals with unobserved places. Rather than marking all runners-up as tied for fourth place, in this case one sets the  $w_i = 0$  for  $i > 3$ . The regression is then not asked to make distinctions between the tied runners-up.

**Breaking Changes**

This package is a work in progress. Expect breaking changes. Please file any bug reports or issues at <https://github.com/shabbychef/ohenery/issues>.

### Legal Mumbo Jumbo

ohenery is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU Lesser General Public License for more details.

### Note

This package is maintained as a hobby.

### Author(s)

Steven E. Pav <shabbychef@gmail.com>

**Maintainer:** Steven E. Pav <shabbychef@gmail.com> ([ORCID](#))

### References

Harville, D. A. "Assigning probabilities to the outcomes of multi-entry competitions." *Journal of the American Statistical Association* 68, no. 342 (1973): 312-316. doi:10.1080/01621459.1973.10482425

Henery, R. J. "Permutation probabilities as models for horse races." *Journal of the Royal Statistical Society: Series B (Methodological)* 43, no. 1 (1981): 86-91. doi:10.1111/j.25176161.1981.tb01153.x

### See Also

Useful links:

- <https://github.com/shabbychef/ohenery>
- Report bugs at <https://github.com/shabbychef/ohenery/issues>

---

ohenery-NEWS

*News for package 'ohenery':*

---

### Description

News for package 'ohenery'

#### **ohenery** Version 0.1.4 (2025-09-07)

- Adding regularization. #9
- Adding hensmfit experts function. #10
- Adding broom tidy and glance methods. #11

#### **ohenery** Version 0.1.3 (2025-04-19)

- Fix issue with factor levels in predict. #8

**ohenery Version 0.1.2 (2024-10-25)**

- Allow warm start in harsm and hensm.

**ohenery Version 0.1.1 (2018-10-14)**

- Change default in harsm and hensm to use unnormalized weights, correcting inference when not all finishes are observed.

**ohenery Initial Version 0.1.0 (2018-10-01)**

- first CRAN release.

---

 race\_data

*Horse Race Data*


---

**Description**

Three weeks of horse race data from tracks worldwide.

**Usage**

```
data(race_data)
```

**Format**

A data.frame object with 36,418 observations and 19 columns.

The columns are defined as follows:

**EventId** An integer ID denoting the event (race). These range from 1 to 4486.

**TrackId** An integer ID number of the the track. There are 64 different tracks represented.

**Type** The type of event, one of “Thoroughbred” or “Harness”.

**RaceNum** The integer race number within a group of races at a track on a given date.

**CorrectedPostTime** The ‘corrected’ post time of the race, in the form %Y-%m-%d %H:%M:%S, presumably in the PDT time zone. Has values like “2019-03-05 02:30:00”.

**Yards** The length of the race, in yards.

**SurfaceText** A string, one of “Turf”, “Dirt”, “All-Weather” or NA.

**HorseName** The string name of the horse.

**HorseId** A unique integer ID for each horse. As different horses can have the same name, this ID is constructed from the name of the Horse, the Sire and the Dam.

**Age** The age of the horse, in integer years, at the time of the event. Typically less than 10.

**Sex** A single character denoting the sex of the horse. I believe the codes are “M” for “Mare” (female four years or older), “G” for “Gelding”, “F” for “Filly” (female under four years of age), “C” for “Colt” (male under four years of age), “H” for “Horse” (male four years of age and up), “R” for “Rig” (hard to explain), “A” for “???”. There are some NA values as well.

**Weight\_lbs** The weight in integer pounds of the jockey and any equipment. Typically around 120.

**PostPosition** The integer starting position of the horse. Typically there is a slight advantage to starting at the first or second post position.

**Medication** One of several codes indicating any medication the horse may be taking at the time of the race. I believe “L” stands for “Lasix”, a common medication for lung conditions that is thought to give horses a slight boost in speed.

**MorningLine** A double indicating the “morning betting line” for win bets on the horse. It is not clear how to interpret this value, perhaps it is return on a dollar. Values range from 0.40 to 80.

**WN\_pool** The total combined pool in win bets, in dollars, on this horse at post time.

**PL\_pool** The total combined pool in place bets, in dollars, on this horse at post time.

**SH\_pool** The total combined pool in show bets, in dollars, on this horse at post time.

**Finish** The integer finishing position of the horse. A 1 means first place. We only collect values of 1, 2, and 3, while the remaining finishing places are unknown and left as NA.

### Note

The author makes no guarantees regarding correctness of this data.

### Author(s)

Steven E. Pav <shabbychef@gmail.com>

### Source

Data were sourced from the web. Don’t ask.

### Examples

```
library(dplyr)
data(race_data)

# compute win bet efficiency
efficiency <- race_data %>%
  group_by(EventId) %>%
  mutate(ImpliedOdds=WN_pool / sum(WN_pool,na.rm=TRUE)) %>%
  ungroup() %>%
  mutate(OddsBucket=cut(ImpliedOdds,c(0,0.05,seq(0.1,1,by=0.10)),include.lowest=TRUE)) %>%
  group_by(OddsBucket) %>%
  summarize(PropWin=mean(as.numeric(coalesce(Finish==1,FALSE)),na.rm=TRUE),
            MedImpl=median(ImpliedOdds,na.rm=TRUE),
            nObs=n()) %>%
  ungroup()

if (require('ggplot2') && require('scales')) {
  efficiency %>%
    ggplot(aes(MedImpl,PropWin,size=nObs)) +
    geom_point() +
    scale_x_sqrt(labels=percent) +
```

```

scale_y_sqrt(labels=percent) +
geom_abline(slope=1,intercept=0,linetype=2,alpha=0.6) +
labs(title='actual win probability versus implied win probability',
      size='# horses',
      x='implied win probability',
      y='observed win probability')
}

```

---

rhenery

*Random generation under the Henery (or Harville) softmax model.*


---

### Description

Given base probabilities, and Henery gamma coefficients, performs random generation, using R's built in rand seed, of the final outcome of a race for each participant.

### Usage

```
rhenery(mu, gamma = NULL)
```

### Arguments

mu	a vector of the probabilities of taking first place.
gamma	a vector of the gamma coefficients. Should have length one less than mu, but if longer the unused elements are ignored. If shorter, we reserve the right to either throw an error or extend out the last gamma element. If not given, the coefficients are assumed to be all one, which is the Harville model.

### Details

Given vectors  $\mu$  and  $\gamma$ , first computes

$$\pi_{1,i} = \frac{\mu_i^{\gamma_1}}{\sum_j \mu_j^{\gamma_1}},$$

then assigns a 1 to participant  $i$  with probability  $\pi_{1,i}$ . The 'winning' participant is then removed from consideration, and the process is repeated using the remaining  $\mu$  and  $\gamma$  vectors.

Typically one has that  $\mu_i = \exp \eta_i$ , for some 'odds',  $\eta_i$ .

When the  $\gamma$  are all one, you recover the Harville softmax model.

### Value

A vector, of the same length as the probabilities, giving the entry of each horse. Note that the expected value of this returned thing makes sense, it is *not* the finished rank ordering of a race.

**Author(s)**

Steven E. Pav <shabbychef@gmail.com>

**See Also**

[rsm](#)

---

rsm

*Generate variates from a softmax distribution.*

---

**Description**

Generate variates from a softmax distribution under Harville or Henery models.

**Usage**

```
rsm(eta, g = NULL, mu = NULL, gamma = NULL)
```

**Arguments**

eta	a vector of the odds. Must be the same length as g if g is given.
g	a vector giving the group indices. If NULL, then we assume only one group is in consideration.
mu	a vector of the probabilities. Must be the same length as g if g is given. If mu and eta are both given, we ignore eta and use mu.
gamma	a vector of the gamma parameters for the Henery model. Typically the first element should be 1. Omit or set all values to 1 to recover the Harville model. The last element will be extended if the gamma is not long enough for a given group. Note that gamma is expected to be very short, and is not 'aligned' with eta in any way.

**Details**

Given the  $\eta$  in odds space, and a grouping variable, returns values in one to the number of elements in a group. That is, we have a permutation of 1 through  $n_g$  on each group as output.

For a single group, the probability that the  $i$ th element of the output is a 1 is equal to

$$\pi_{1,i} = \frac{\mu_i^{\gamma_1}}{\sum_j \mu_j^{\gamma_1}}.$$

Once an element has been selected to have output 1, remove it from the set and iterate, but with the next  $\gamma$  elements.

**Value**

a vector of integers, each a permutation of one through the number of elements in each group.

**Note**

The output of this function satisfies a kind of order invariance that [rhenery](#) does not, at the cost of some computational inefficiency. Namely that for a fixed randseed, and for distinct eta, the output is equivariant with respect to permutation of the vector eta when there is only one group.

Regarding the ‘direction’, we associate higher odds with a smaller outcome. That is, the ith element of the output encodes the place that the ith participant took in the simulated ‘race’; it should be small if the odds for that participant are very high.

**Author(s)**

Steven E. Pav <shabbychef@gmail.com>

**Examples**

```
# simple use
set.seed(1234)
g <- ceiling(seq(1,10,by=0.1))
eta <- rnorm(length(g))
y <- rsm(eta,g=g)

# same same:
set.seed(235)
y1 <- rsm(eta,g=g)
set.seed(235)
y2 <- rsm(g=g,mu=smax(eta,g=g))
y1 - y2

# the default model is Harville
set.seed(1212)
y1 <- rsm(eta,g=g)
set.seed(1212)
y2 <- rsm(eta,g=g,gamma=c(1,1,1,1))
y1 - y2

# repeat several times with the cards stack against
# early runners
set.seed(1234)
colMeans(t(replicate(1000,rsm(sort(rnorm(10,sd=1.5)),g=rep(1,10))))))

# illustrate the invariance
mu <- (1:10) / 55
set.seed(1414)
y1 <- rsm(mu=mu,gamma=c(1,1,1))
set.seed(1414)
y2 <- rev(rsm(mu=rev(mu),gamma=c(1,1,1)))
y1 - y2

nfeat <- 5
set.seed(1234)
g <- ceiling(seq(0.1,1000,by=0.1))
```

```

X <- matrix(rnorm(length(g) * nfeat),ncol=nfeat)
beta <- rnorm(nfeat)
eta <- X %*% beta
y <- rsm(eta,g=g)

idx <- order(g,y,decreasing=TRUE) - 1
foeey <- harsmlik(g,idx,eta,deleta=X)
set.seed(3493)
dib <- rnorm(length(beta))
xv1 <- seq(-0.01,0.01,length.out=301)
rsu <- sapply(xv1,
              function(del) {
                beta1 <- beta + del * dib
                eta1 <- X %*% beta1
                foeey2 <- harsmlik(g,idx,eta1,deleta=X)
                as.numeric(foeey2) - as.numeric(foeey)
              })
drv <- sapply(xv1,
              function(del) {
                beta1 <- beta + del * dib
                eta1 <- X %*% beta1
                foeey2 <- harsmlik(g,idx,eta1,deleta=X)
                sum(attr(foeey2,'gradient') * dib)
              })

if (require('ggplot2') && require('dplyr')) {
  bestx <- xv1[which.max(rsu)]
  ph <- data.frame(x=xv1,lik=rsu,grd=drv) %>%
    ggplot(aes(x=x,y=lik)) +
    geom_point() +
    geom_line(aes(y=grd/200)) +
    geom_vline(xintercept=bestx,linetype=2,alpha=0.5) +
    geom_hline(yintercept=0,linetype=2,alpha=0.5)
  print(ph)
}

if (require('dplyr') && require('knitr')) {
  # expect this to be very small, almost always 1
  set.seed(1234)
  simdraw <- replicate(10000,{
    rsm(eta=c(100,rnorm(7)))[1]
  })

  as.data.frame(table(simdraw)) %>%
    mutate(prob=Freq / sum(Freq)) %>%
    knitr::kable()

  # expect this to be uniform on 2 through 8
  set.seed(1234)
  simdraw <- replicate(10000,{
    rsm(eta=c(100,rnorm(7)))[2]
  })
}

```

```

as.data.frame(table(simdraw)) %>%
  mutate(prob=Freq / sum(Freq)) %>%
  knitr::kable()
}

```

---

smax

*The softmax function.*


---

### Description

The softmax function: exponentiate a vector and then normalize.

### Usage

```
smax(eta, g = NULL)
```

### Arguments

**eta** numeric array of the odds. The odds are de-meant within each group.

**g** a vector giving the group indices. If NULL, then we assume only one group is in consideration.

### Details

Given vector  $\eta$  for a single group, essentially computes vector  $\mu$  defined by

$$\mu_i = \frac{\exp \eta_i}{\sum_j \exp \eta_j}.$$

Note that this computation should be invariant with respect to level shifts of the  $\eta$ , and thus we de-mean the odds first.

### Value

the exponentiated data normalized. For the row-wise version, each row is soft maxed.

### Note

This function can deal with overflow in a semi-coherent way.

### Author(s)

Steven E. Pav <shabbychef@gmail.com>

### See Also

[normalize](#), [inv\\_smax](#).

**Examples**

```
# we can deal with large values:
set.seed(2345)
eta <- rnorm(12,sd=1000)
smax(eta)
```

---

smlk

*Softmax log likelihood under Harville and Henery Models.*


---

**Description**

Compute the softmax log likelihood and gradient of the same.

**Usage**

```
harsmlk(g, idx, eta, wt = NULL, deleta = NULL)
hensmlk(g, idx, eta, gamma, wt = NULL, deleta = NULL)
```

**Arguments**

<code>g</code>	a vector giving the group indices.
<code>idx</code>	a vector of integers, the same length as <code>g</code> , which describes the reverse sort order on the observations, first by group, then by place within the group. That is, the element <code>idx[0]</code> is the index of the last place finisher in the group <code>g[0]</code> ; then <code>idx[1]</code> is the index of the next-to-last place finisher in <code>g[1]</code> (assuming it equals <code>g[0]</code> ), and so on. The <code>idx</code> shall be zero based.
<code>eta</code>	a vector of the odds. Must be the same length as <code>g</code> .
<code>wt</code>	an optional vector of non-negative elements, the same length as <code>g</code> , giving the observation level weights. We then compute a weighted log likelihood, where the weights are in each summand. The weights should probably have mean 1, but that's just, like, my opinion, man. Negative weights throw an error. Note that the weights for the last place in each group have no effect on the computation.
<code>deleta</code>	an optional matrix whose row count equals the number of elements of <code>g</code> , <code>idx</code> , and <code>eta</code> . The rows of <code>deleta</code> are the derivatives of <code>eta</code> with respect to some $z$ . This is used to then maximize likelihood with respect to $z$ .
<code>gamma</code>	a vector of the gamma parameters. It is assumed that the first element is $\gamma_2$ , while the last element is applied to all higher order tie breaks.

**Details**

Given vectors  $g$ ,  $\eta$  and optionally the gradient of  $\eta$  with respect to some coefficients, computes the log likelihood under the softmax. The user must provide a reverse ordering as well, which is sorted first by the groups,  $g$ , and then, within a group, in increasing quality order. For a race, this means that the index is in order from the last place to the first place in that race, where the group numbers correspond to one race.

Under the Harville model, the log likelihood on a given group, where we are indexing in *forward* order, is

$$\left( \eta_1 - \log \sum_{j \geq 1} \mu_j \right) + \left( \eta_2 - \log \sum_{j \geq 2} \mu_j \right) + \dots$$

where  $\mu_i = \exp \eta_i$ . By “forward order”, we mean that  $\eta_1$  corresponds to the participant taking first place within that group,  $\eta_2$  took second place, and so on.

The Henery model log likelihood takes the form

$$\left( \eta_1 - \log \sum_{j \geq 1} \mu_j \right) + \left( \gamma_2 \eta_2 - \log \sum_{j \geq 2} \mu_j^{\gamma_2} \right) + \dots$$

for gamma parameters,  $\gamma$ . The Henery model corresponds to the Harville model where all the gammas equal 1.

Weights in weighted estimation apply to each summand. The weight for the last place participant in a group is irrelevant. The weighted log likelihood under the Harville model is

$$w_1 \left( \eta_1 - \log \sum_{j \geq 1} \mu_j \right) + w_2 \left( \eta_2 - \log \sum_{j \geq 2} \mu_j \right) + \dots$$

One should think of the weights as applying to the outcome, not the participant.

## Value

The log likelihood. If `delta` is given, we add an attribute to the scalar number, called `gradient` giving the derivative. For the Henery model we also include a term of `gradgamma` which is the gradient of the log likelihood with respect to the gamma vector.

## Note

The likelihood function does not yet support ties.

To avoid incorrect inference when only the top performers are recorded, and all others are effectively tied, one should use weighting. Set the weights to zero for participants who are tied non-winners, and one for the rest. So for example, if you observe the Gold, Silver, and Bronze medal winners of an Olympic event that had a starting field of 12 participants, set weights to 1 for the medal winners, and 0 for the others. Note that the weights do not attach to the participants, they attach to the place they took.

## Author(s)

Steven E. Pav <shabbychef@gmail.com>

## References

- Harville, D. A. "Assigning probabilities to the outcomes of multi-entry competitions." *Journal of the American Statistical Association* 68, no. 342 (1973): 312-316. doi:10.1080/01621459.1973.10482425
- Henery, R. J. "Permutation probabilities as models for horse races." *Journal of the Royal Statistical Society: Series B (Methodological)* 43, no. 1 (1981): 86-91. doi:10.1111/j.25176161.1981.tb01153.x

**Examples**

```

# a garbage example
set.seed(12345)
g <- as.integer(sort(ceiling(20 * runif(100))))
idx <- as.integer(rev(1:length(g)) - 1L)
eta <- rnorm(length(g))
foo <- harsmlik(g=g,idx=idx,eta=eta,deleta=NULL)

# an example with a real idx
nfeat <- 5
set.seed(1234)
g <- ceiling(seq(0.1,1000,by=0.1))
X <- matrix(rnorm(length(g) * nfeat),ncol=nfeat)
beta <- rnorm(nfeat)
eta <- X %*% beta
y <- rsm(eta,g)

idx <- order(g,y,decreasing=TRUE) - 1
foores <- harsmlik(g,idx,eta,deleta=X)

# now reweight them
wt <- runif(length(g))
wt <- wt / mean(wt) # mean 1 is recommended
foores <- harsmlik(g,idx,eta,wt=wt)

# try hensmlik
foores <- hensmlik(g,idx,eta,gamma=c(0.9,0.8,1),wt=wt)

# check the value of the gradient by numerical approximation

nfeat <- 8
set.seed(321)
g <- ceiling(seq(0.1,1000,by=0.1))
X <- matrix(rnorm(length(g) * nfeat),ncol=nfeat)
beta <- rnorm(nfeat)
eta <- X %*% beta
y <- rsm(eta,g)

idx <- order(g,y,decreasing=TRUE) - 1
if (require(numDeriv)) {

  fastval <- attr(harsmlik(g,idx,eta,deleta=X),'gradient')
  numap <- grad(function(beta,g,idx,X) {
    eta <- X %*% beta
    as.numeric(harsmlik(g,idx,eta))
  },
  x=beta,g=g,idx=idx,X=X)
  rbind(fastval,numap)
}

```

---

`tidy.harsm`*Tidy*

---

## Description

Tidy a harsm, hensm or linodds model.

## Usage

```
## S3 method for class 'harsm'  
tidy(x, ...)
```

```
## S3 method for class 'hensm'  
tidy(x, ...)
```

## Arguments

`x` an object of type harsm or hensm  
`...` arguments for generic consistency.

## Details

Returns a table with information on the fit coefficients and standard errors of an estimated Harville or Henery model.

## Value

A tidy `tibble::tibble()` with fields

**term** The name of the estimated parameter. Betas typically come before gammas.

**estimate** The estimated parameter.

**std.error** The standard error of the estimate.

**statistic** The z-statistic of the estimate.

## Author(s)

Steven E. Pav <shabbychef@gmail.com>

## See Also

[tidy.maxLik](#).

**Examples**

```
# softmax on the Best Picture data
data(best_picture)
df <- best_picture
df$place <- ifelse(df$winner,1,2)
df$weight <- ifelse(df$winner,1,0)

fmla <- place ~ nominated_for_BestDirector + nominated_for_BestActor + Drama
fit0 <- harsm(fmla,data=df,group=year,weights=weight)
print(tidy(fit0))
```

# Index

- \* **data**
  - best\_picture, 4
  - diving, 6
  - race\_data, 29
- \* **fitting**
  - harsm, 10
  - harsmfit, 14
  - hensm, 18
  - hensmfit, 22
- \* **package**
  - ohenery, 27
- \* **probability**
  - rhensy, 31
  - rsm, 32
- as.data.frame, 11, 19
- as.linodds, 2
- best\_picture, 4
- coef.linodds (as.linodds), 2
- diving, 6
- erank, 7, 17, 18, 24
- formula, 3, 10, 19
- glance.harsm, 9
- glance.hensm (glance.harsm), 9
- glance.maxLik, 9
- harsm, 3, 10, 20, 21
- harsm\_invlink, 3, 17
- harsmfit, 13, 14
- harsmlik, 13, 16, 17
- harsmlik (smlk), 36
- hensm, 3, 18
- hensmfit, 22
- hensmlik, 24
- hensmlik (smlk), 36
- inv\_smax, 25, 35
- maxBHHH, 16, 24
- maxLik, 12, 16, 17, 20, 24
- normalize, 26, 35
- ohenery, 27
- ohenery-NEWS, 28
- ohenery-package (ohenery), 27
- predict.linodds (as.linodds), 2
- print.harsm (harsm), 10
- print.hensm (hensm), 18
- race\_data, 29
- rhensy, 31, 33
- rsm, 32, 32
- smax, 3, 26, 35
- smlk, 21, 36
- tidy.harsm, 39
- tidy.hensm (tidy.harsm), 39
- tidy.maxLik, 39
- vcov.harsm (harsm), 10
- vcov.hensm (hensm), 18