

# Package ‘updog’

September 10, 2019

**Title** Flexible Genotyping for Polyploids

**Version** 1.1.1

**Description** Implements empirical Bayes approaches to genotype polyploids from next generation sequencing data while accounting for allelic bias, overdispersion, and sequencing error. The main function is `flexdog()`, which allows the specification of many different genotype distributions. An experimental function that takes into account varying levels of relatedness is implemented in `mupdog()`. Also provided are functions to simulate genotypes, `rgeno()`, and read-counts, `rflexdog()`, as well as functions to calculate oracle genotyping error rates, `oracle_mis()`, and correlation with the true genotypes, `oracle_cor()`. These latter two functions are useful for read depth calculations. Run `browseVignettes(package = "updog")` in R for example usage. See Gerard et al. (2018) <doi:10.1534/genetics.118.301468> and Gerard and Ferrao (2019) <doi:10.1101/751784> for details on the implemented methods.

**Depends** R (>= 3.4.0)

**License** GPL-3

**BugReports** <http://github.com/dcgerard/updog/issues>

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 6.1.1

**LinkingTo** Rcpp, RcppArmadillo

**Imports** Rcpp (>= 0.12.16), RcppArmadillo, assertthat, foreach, doParallel, ggplot2, ggthemes, stringr

**Suggests** covr, testthat, SuppDists, Rmpfr, CVXR, knitr, rmarkdown, tidyverse

**VignetteBuilder** knitr

**NeedsCompilation** yes

**Author** David Gerard [aut, cre] (<<https://orcid.org/0000-0001-9450-5023>>)

**Maintainer** David Gerard <gerard.1787@gmail.com>

**Repository** CRAN

**Date/Publication** 2019-09-09 23:30:02 UTC

## R topics documented:

ashpen_fun . . . . .	4
compute_all_log_bb . . . . .	5
compute_all_phifk . . . . .	5
compute_all_post_prob . . . . .	6
convolve_up . . . . .	7
dbernbinom . . . . .	8
dbetabinom . . . . .	8
dbetabinom_alpha_beta_double . . . . .	10
dbetabinom_double . . . . .	11
dc_dtau . . . . .	11
df_deps . . . . .	12
dlbeta_dc . . . . .	12
dlbeta_deps . . . . .	13
dlbeta_dh . . . . .	14
dlbeta_dtau . . . . .	14
dlbeta_dxi . . . . .	15
doutdist . . . . .	16
dpen_deps . . . . .	16
dpen_dh . . . . .	17
dr_pen . . . . .	18
dxi_df . . . . .	18
dxi_dh . . . . .	19
elbo . . . . .	19
eta_double . . . . .	20
eta_fun . . . . .	21
expit . . . . .	21
f1_obj . . . . .	22
flexdog . . . . .	22
flexdog_full . . . . .	27
flexdog_obj . . . . .	33
flexdog_obj_out . . . . .	34
flex_update_pivec . . . . .	35
get_bivalent_probs . . . . .	35
get_bivalent_probs_dr . . . . .	36
get_conv_inner_weights . . . . .	37
get_dimname . . . . .	38
get_hyper_weights . . . . .	38
get_inner_weights . . . . .	39
get_probk_vec . . . . .	40
get_q_array . . . . .	40
get_uni_rep . . . . .	41

get_wik_mat . . . . .	42
get_wik_mat_out . . . . .	43
grad_for_eps . . . . .	44
grad_for_mu_sigma2 . . . . .	45
grad_for_mu_sigma2_wrapper . . . . .	45
grad_for_weighted_lbb . . . . .	46
grad_for_weighted_lnorm . . . . .	47
initialize_pivec . . . . .	47
is.flexdog . . . . .	48
is.mupdog . . . . .	49
logit . . . . .	49
log_sum_exp . . . . .	50
log_sum_exp_2 . . . . .	50
mupdog . . . . .	51
mupout . . . . .	54
obj_for_alpha . . . . .	55
obj_for_eps . . . . .	56
obj_for_mu_sigma2 . . . . .	57
obj_for_mu_sigma2_wrapper . . . . .	57
obj_for_rho . . . . .	58
obj_for_weighted_lbb . . . . .	59
obj_for_weighted_lnorm . . . . .	59
oracle_cor . . . . .	60
oracle_cor_from_joint . . . . .	61
oracle_joint . . . . .	62
oracle_mis . . . . .	64
oracle_mis_from_joint . . . . .	65
oracle_mis_vec . . . . .	67
oracle_mis_vec_from_joint . . . . .	68
oracle_plot . . . . .	69
pbetabinom_double . . . . .	70
pen_bias . . . . .	71
pen_seq_error . . . . .	71
pivec_from_segmat . . . . .	72
plot.flexdog . . . . .	73
plot.mupdog . . . . .	74
plot_geno . . . . .	75
post_prob . . . . .	76
pp_brent_obj . . . . .	77
qbetabinom_double . . . . .	78
rbetabinom_int . . . . .	78
rflexdog . . . . .	79
rgeno . . . . .	80
snpdat . . . . .	82
summary.mupdog . . . . .	83
uitdewilligen . . . . .	84
uni_em . . . . .	85
uni_em_const . . . . .	86

uni_obj	87
uni_obj_const	87
update_dr	88
update_pp_f1	89
update_pp_s1	90
update_R	91
updog	91
wem	93
xi_double	94
xi_fun	95

<b>Index</b>	<b>96</b>
--------------	-----------

---

ashpen_fun	<i>Penalty on pivec used when model = "ash" in flexdog.</i>
------------	---

---

### Description

Penalty on pivec used when model = "ash" in flexdog.

### Usage

```
ashpen_fun(lambda, pivec)
```

### Arguments

lambda	The penalty.
pivec	The vector of mixing proportions for the component discrete uniform distributions.

### Value

A penalty on the ash mixing weights.

### Author(s)

David Gerard

---

compute\_all\_log\_bb      *Calculates the log-density for every individual by snp by dosage level.*

---

### Description

Calculates the log-density for every individual by snp by dosage level.

### Usage

```
compute_all_log_bb(refmat, sizemat, ploidy, seq, bias, od)
```

### Arguments

refmat	A matrix of reference counts. The rows index the individuals and the columns index the SNPs.
sizemat	A matrix of total counts. The rows index the individuals and the columns index the SNPs. Should have the same dimensions as refmat.
ploidy	The ploidy of the species. To estimate the ploidy, re-run mupdog at various ploidy levels and choose the one with the largest ELBO. This assumes that the ploidy is the same for all individuals in the sample.
seq	A vector of initial sequencing errors. Should be the same length as the number of columns of refmat (number of SNPs). Must be between 0 and 1.
bias	A vector of initial bias parameters. Should be the same length as the number of columns of refmat (number of SNPs). Must be greater than 0.
od	A vector of initial overdispersion parameters. Should be the same length as the number of columns of refmat (number of SNPs). Must be between 0 and 1.

### Value

A three dimensional array. The rows index the individuals, the columns index the SNPs, and the third dimension indexes the genotypes. This is the log-likelihood for each individual/snp/genotype combination.

---

compute\_all\_phifk      *Computes*

$$\Phi^{-1}(F(k|K, \alpha_j, \rho_i))$$

*for all possible (i,j,k).*

---

### Description

Computes

$$\Phi^{-1}(F(k|K, \alpha_j, \rho_i))$$

for all possible (i,j,k).

**Usage**

```
compute_all_phifk(alpha, rho, ploidy)
```

**Arguments**

alpha	A vector whose jth element is the allele frequency of SNP j.
rho	A vector whose ith element is the inbreeding coefficient of individual i.
ploidy	The ploidy of the species.

**Value**

A three dimensional array. The rows index the individuals, the columns index the SNPs, and the third dimension indexes the genotypes. Computes the "continuous genotype".

**Author(s)**

David Gerard

---

compute\_all\_post\_prob *Computes every posterior probability for each dosage level for each individual at each SNP.*

---

**Description**

Computes every posterior probability for each dosage level for each individual at each SNP.

**Usage**

```
compute_all_post_prob(ploidy, mu, sigma2, alpha, rho)
```

**Arguments**

ploidy	The ploidy of the species.
mu	A matrix of variational posterior means. The rows index the individuals and the columns index the SNPs.
sigma2	A matrix of variational posterior variances. The rows index the individuals and the columns index the SNPs.
alpha	A vector of allele frequencies for all SNPs.
rho	A vector of inbreeding coefficients for all individuals.

**Value**

An array. The rows index the individuals, the columns index the SNPS, and the third dimension indexes the genotypes. Element (i, j, k) is the return of [post\\_prob](#).

**Author(s)**

David Gerard

---

convolve_up	<i>Convolution between two discrete probability mass functions with support on 0:K.</i>
-------------	---

---

**Description**

Convolution between two discrete probability mass functions with support on 0:K.

**Usage**

```
convolve_up(x, y)
```

**Arguments**

x	The first probability vector. The ith element is the probability of i - 1.
y	The second probability vector. The ith element is the probability of i - 1.

**Value**

A vector that is the convolution of x and y. The ith element is the probability of i - 1.

**Author(s)**

David Gerard

**See Also**

[convolve](#) for a more generic convolution function.

**Examples**

```
x <- c(1 / 6, 2 / 6, 3 / 6)
y <- c(1 / 9, 2 / 9, 6 / 9)
convolve_up(x, y)
stats::convolve(x, rev(y), type = "o")
```

---

dbernbinom	<i>Special case of betabinomial where the beta is bernoulli mu.</i>
------------	---

---

**Description**

Special case of betabinomial where the beta is bernoulli mu.

**Usage**

```
dbernbinom(x, size, mu, log)
```

**Arguments**

x	The quantile.
size	The total number of draws.
mu	The mean of the beta.
log	A logical. Should we return the log of the density TRUE or not FALSE?

**Value**

The density of the Bernoulli-binomial.

**Author(s)**

David Gerard

---

dbetabinom	<i>The Beta-Binomial Distribution</i>
------------	---------------------------------------

---

**Description**

Density, distribution function, quantile function and random generation for the beta-binomial distribution when parameterized by the mean  $\mu$  and the overdispersion parameter  $\rho$  rather than the typical shape parameters.

**Usage**

```
dbetabinom(x, size, mu, rho, log)
pbetabinom(q, size, mu, rho, log_p)
qbetabinom(p, size, mu, rho)
rbetabinom(n, size, mu, rho)
```



**Arguments**

<code>x, q</code>	A vector of quantiles.
<code>size</code>	A vector of sizes.
<code>mu</code>	Either a scalar of the mean for each observation, or a vector of means of each observation, and thus the same length as <code>x</code> and <code>size</code> . This must be between 0 and 1.
<code>rho</code>	Either a scalar of the overdispersion parameter for each observation, or a vector of overdispersion parameters of each observation, and thus the same length as <code>x</code> and <code>size</code> . This must be between 0 and 1.
<code>log, log_p</code>	A logical vector either of length 1 or the same length as <code>x</code> and <code>size</code> . This determines whether to return the log probabilities for all observations (in the case that its length is 1) or for each observation (in the case that its length is that of <code>x</code> and <code>size</code> ).
<code>p</code>	A vector of probabilities.
<code>n</code>	The number of observations.

**Details**

Let  $\mu$  and  $\rho$  be the mean and overdispersion parameters. Let  $\alpha$  and  $\beta$  be the usual shape parameters of a beta distribution. Then we have the relation

$$\mu = \alpha / (\alpha + \beta),$$

and

$$\rho = 1 / (1 + \alpha + \beta).$$

This necessarily means that

$$\alpha = \mu(1 - \rho) / \rho,$$

and

$$\beta = (1 - \mu)(1 - \rho) / \rho.$$

**Value**

Either a random sample (`rbetabinom`), the density (`dbetabinom`), the tail probability (`pbetabinom`), or the quantile (`qbetabinom`) of the beta-binomial distribution.

**Functions**

- `dbetabinom`: Density function.
- `pbetabinom`: Distribution function.
- `qbetabinom`: Quantile function.
- `rbetabinom`: Random generation.

**Author(s)**

David Gerard

**Examples**

```
x <- rbetabinom(n = 10, size = 10, mu = 0.1, rho = 0.01)
dbetabinom(x = 1, size = 10, mu = 0.1, rho = 0.01, log = FALSE)
pbetabinom(q = 1, size = 10, mu = 0.1, rho = 0.01, log_p = FALSE)
qbetabinom(p = 0.6, size = 10, mu = 0.1, rho = 0.01)
```

---

dbetabinom\_alpha\_beta\_double

*Density function of betabinomial with the shape parameterizations*

---

**Description**

Density function of betabinomial with the shape parameterizations

**Usage**

```
dbetabinom_alpha_beta_double(x, size, alpha, beta, log)
```

**Arguments**

x	The quantile.
size	The total number of draws.
alpha	The first shape parameter.
beta	The second shape parameter.
log	A logical. Should we return the log of the density TRUE or not FALSE?

**Value**

The density of the beta-binomial.

**Author(s)**

David Gerard

---

dbetabinom\_double      *The density function of the beta-binomial distribution.*

---

**Description**

The density function of the beta-binomial distribution.

**Usage**

```
dbetabinom_double(x, size, mu, rho, log)
```

**Arguments**

x	The quantile.
size	The total number of draws.
mu	The mean of the beta.
rho	The overdispersion parameter of the beta.
log	A logical. Should we return the log of the density TRUE or not FALSE?

**Value**

The density of the beta-binomial.

**Author(s)**

David Gerard

---

dc\_dtau      *Derivative of  $c = (1 - \tau)/\tau$  with respect to  $\tau$ .*

---

**Description**

Derivative of  $c = (1 - \tau)/\tau$  with respect to  $\tau$ .

**Usage**

```
dc_dtau(tau)
```

**Arguments**

tau	The overdispersion parameter.
-----	-------------------------------

**Value**

A double.

**Author(s)**

David Gerard

**See Also**[dlbeta\\_dc](#), [dlbeta\\_dtau](#)

---

`df_deps`*Derivative of f with respect to eps.*

---

**Description**

Derivative of f with respect to eps.

**Usage**`df_deps(p, eps)`**Arguments**

`p`                    The allele dosage.  
`eps`                    The sequencing error rate.

**Value**

A double.

**Author(s)**

David Gerard

---

`dlbeta_dc`*Derivative of the log-beta density with respect to c where  $c = (1 - \tau)/\tau$  where  $\tau$  is the overdispersion parameter.*

---

**Description**Derivative of the log-beta density with respect to c where  $c = (1 - \tau)/\tau$  where  $\tau$  is the overdispersion parameter.**Usage**`dlbeta_dc(x, n, xi, c)`

**Arguments**

x	The number of successes observed
n	The total number of trials observed.
$\bar{x}$	The mean of the beta-binomial.
c	$(1 - \tau)/\tau$ where $\tau$ is the overdispersion parameter.

**Value**

A double.

**Author(s)**

David Gerard

**See Also**

[dbetabinom\\_double](#), [dlbeta\\_dtau](#), [dc\\_dtau](#).

---

dlbeta_deps	<i>Derivative of the log-beta-binomial density with respect to the sequencing error rate.</i>
-------------	---

---

**Description**

Derivative of the log-beta-binomial density with respect to the sequencing error rate.

**Usage**

```
dlbeta_deps(x, n, p, eps, h, tau)
```

**Arguments**

x	The number of successes.
n	The number of trials.
p	The allele dosage.
eps	The sequencing error rate
h	The bias parameter.
tau	The overdispersion parameter.

**Value**

A double.

**Author(s)**

David Gerard

---

dlbeta_dh	<i>Derivative of log-betabinomial density with respect to bias parameter.</i>
-----------	---

---

**Description**

Derivative of log-betabinomial density with respect to bias parameter.

**Usage**

dlbeta\_dh(x, n, p, eps, h, tau)

**Arguments**

x	The number of successes.
n	The number of trials.
p	The allele dosage.
eps	The sequencing error rate
h	The bias parameter.
tau	The overdispersion parameter.

**Value**

A double.

**Author(s)**

David Gerard

---

dlbeta_dtau	<i>Derivative of the log-beta-binomial density with respect to the overdispersion parameter.</i>
-------------	--

---

**Description**

Derivative of the log-beta-binomial density with respect to the overdispersion parameter.

**Usage**

dlbeta\_dtau(x, n, p, eps, h, tau)

**Arguments**

x	The number of successes.
n	The number of trials.
p	The allele dosage.
eps	The sequencing error rate
h	The bias parameter.
tau	The overdispersion parameter.

**Value**

A double.

**Author(s)**

David Gerard

**See Also**

[dlbeta\\_dc](#), [dc\\_dtau](#), [dbetabinom\\_double](#).

---

dlbeta_dxi	<i>Derivative of the log-betabinomial density with respect to the mean of the underlying beta.</i>
------------	--

---

**Description**

Derivative of the log-betabinomial density with respect to the mean of the underlying beta.

**Usage**

```
dlbeta_dxi(x, n, xi, tau)
```

**Arguments**

x	The number of successes.
n	The number of trials.
xi	The mean of the underlying beta.
tau	The overdispersion parameter.

**Value**

A double.

**Author(s)**

David Gerard

---

doutdist	<i>The outlier distribution we use. Right now it is just a beta binomial with mean 1/2 and od 1/3 (so underlying beta is just a uniform from 0 to 1).</i>
----------	---

---

**Description**

The outlier distribution we use. Right now it is just a beta binomial with mean 1/2 and od 1/3 (so underlying beta is just a uniform from 0 to 1).

**Usage**

doutdist(x, n, logp)

**Arguments**

x	The number of reference counts.
n	The total number of read-counts.
logp	Return the log density TRUE or not FALSE?

**Value**

A double. The outlier density value.

**Author(s)**

David Gerard

---

dpen_deps	<i>Derivative of</i>
	$-\log(\epsilon(1 - \epsilon)) - (\text{logit}(\epsilon) - \mu_\epsilon)^2 / (2\sigma_\epsilon^2)$
	<i>with respect to <math>\epsilon</math>.</i>

---

**Description**

Derivative of

$$-\log(\epsilon(1 - \epsilon)) - (\text{logit}(\epsilon) - \mu_\epsilon)^2 / (2\sigma_\epsilon^2)$$

with respect to  $\epsilon$ .

**Usage**

dpen\_deps(eps, mu\_eps, sigma2\_eps)



**Arguments**

eps	The current sequencing error rate.
mu_eps	The mean of the logit of the sequencing error rate.
sigma2_eps	The variance of the logit of the sequencing error rate.

**Value**

A double.

**Author(s)**

David Gerard

**See Also**

[pen\\_seq\\_error](#) which this is a derivative for.

---

dpen\_dh

*Derivative of*

$$-\log(h) - (\log(h) - \mu_h)^2 / (2\sigma_h^2)$$

*with respect to h.*

---

**Description**

Derivative of

$$-\log(h) - (\log(h) - \mu_h)^2 / (2\sigma_h^2)$$

with respect to  $h$ .

**Usage**

dpen\_dh(h, mu\_h, sigma2\_h)

**Arguments**

h	The current bias parameter.
mu_h	The mean of the log-bias.
sigma2_h	The variance of the log-bias.

**Value**

A double.

**Author(s)**

David Gerard

**See Also**

[pen\\_bias](#) which this is a derivative for.

---

dr_pen	<i>Penalty used in <a href="#">update_dr</a>.</i>
--------	---

---

**Description**

A dirichlet prior on pairweights. Returns log density.

**Usage**

```
dr_pen(pairweights, mixing_pen)
```

**Arguments**

pairweights	The mixing proportions to penalize.
mixing_pen	The corresponding penalties.

**Author(s)**

David Gerard

**See Also**

[update\\_dr](#)

---

dxi_df	<i>Derivative of xi with respect to f.</i>
--------	--

---

**Description**

Derivative of xi with respect to f.

**Usage**

```
dxi_df(h, f)
```

**Arguments**

h	The bias parameter.
f	The post-sequencing error rate adjusted probability of an A.

**Value**

A double.

**Author(s)**

David Gerard

---

dxi_dh	<i>Derivative of xi-function with respect to bias parameter.</i>
--------	--

---

**Description**

Derivative of xi-function with respect to bias parameter.

**Usage**

dxi\_dh(p, eps, h)

**Arguments**

p	The dosage (between 0 and 1).
eps	The sequencing error rate.
h	The bias parameter.

**Value**

A double.

**Author(s)**

David Gerard

---

elbo	<i>The evidence lower bound</i>
------	---------------------------------

---

**Description**

The evidence lower bound

**Usage**

```
elbo(warray, lbeta_array, cor_inv, postmean, postvar, bias, seq, mean_bias,
     var_bias, mean_seq, var_seq, ploidy)
```

**Arguments**

warray	An three-way array. The (i,j,k)th entry is the variational posterior probability that individual i at SNP j has dosage k - 1. See <a href="#">compute_all_post_prob</a> .
lbeta_array	A three-way array. The (i,j,k)th entry is the log-density of the betabinomial for individual i at SNP j and dosage k - 1. See <a href="#">compute_all_log_bb</a> .
cor_inv	The inverse of the correlation matrix.
postmean	A matrix. The (i,j)th entry is the variational posterior mean for individual i at SNP j.
postvar	A matrix. The (i,j)th entry is the variational posterior variance for individual i at SNP j.
bias	A vector. The jth entry is the allele bias for SNP j.
seq	A vector. The jth entry is the sequencing error rate at SNP j.
mean_bias	The prior mean on the log-bias.
var_bias	The prior variance on the log-bias.
mean_seq	The prior mean on the logit of the sequencing error rate.
var_seq	The prior variance on the logit of the sequencing error rate.
ploidy	The ploidy of the species.

**Value**

A double. The evidence lower-bound that [mupdog](#) maximizes.

**Author(s)**

David Gerard

---

eta\_double                      *Adjusts allele dosage p by the sequencing error rate eps.*

---

**Description**

Adjusts allele dosage p by the sequencing error rate eps.

**Usage**

```
eta_double(p, eps)
```

**Arguments**

p	The allele dosage.
eps	The sequencing error rate.

**Value**

The probability of a reference read adjusted by the sequencing error rate.

**Author(s)**

David Gerard

---

eta_fun	<i>Adjusts allele dosage p by the sequencing error rate eps.</i>
---------	--

---

**Description**

Adjusts allele dosage p by the sequencing error rate eps.

**Usage**

```
eta_fun(p, eps)
```

**Arguments**

p	A vector of allele dosages.
eps	The sequencing error rate. Must either be of length 1 or the same length as p.

**Value**

A vector of probabilities of a reference read adjusted by the sequencing error rate.

**Author(s)**

David Gerard

---

expit	<i>The expit (logistic) function.</i>
-------	---------------------------------------

---

**Description**

The expit (logistic) function.

**Usage**

```
expit(x)
```

**Arguments**

x	A double.
---	-----------

**Value**

The expit (logistic) of  $x$ .

**Author(s)**

David Gerard

---

f1_obj	<i>Objective for mixture of known dist and uniform dist.</i>
--------	--

---

**Description**

Objective for mixture of known dist and uniform dist.

**Usage**

```
f1_obj(alpha, pvec, weight_vec)
```

**Arguments**

alpha	The mixing weight.
pvec	The known distribution (e.g. from assuming an F1 population).
weight_vec	A vector of weights.

**Value**

The objective when updating pivec when model = "f1" or model = "s1" in [flexdog\\_full](#).

**Author(s)**

David Gerard

---

flexdog	<i>Flexible genotyping for polyploids from next-generation sequencing data.</i>
---------	---

---

**Description**

Genotype polyploid individuals from next generation sequencing (NGS) data while assuming the genotype distribution is one of several forms. flexdog does this while accounting for allele bias, overdispersion, sequencing error, and possibly outlying observations (if model = "f1" or model = "s1"). The method is described in detail in Gerard et. al. (2018) and Gerard and Ferrão (2019).

**Usage**

```
flexdog(refvec, sizevec, ploidy, model = c("norm", "hw", "bb", "ash",
    "s1", "s1pp", "f1", "f1pp", "flex", "uniform", "custom"), p1ref = NULL,
    p1size = NULL, p2ref = NULL, p2size = NULL, snpname = NULL,
    bias_init = exp(c(-1, -0.5, 0, 0.5, 1)), verbose = TRUE,
    outliers = FALSE, prior_vec = NULL, ...)
```

**Arguments**

refvec	A vector of counts of reads of the reference allele.
sizevec	A vector of total counts.
ploidy	The ploidy of the species. Assumed to be the same for each individual.
model	What form should the prior (genotype distribution) take? See Details for possible values.
p1ref	The reference counts for the first parent if model = "f1" (or model = "f1pp"), or for the only parent if model = "s1" (or model = "s1pp").
p1size	The total counts for the first parent if model = "f1" (or model = "f1pp"), or for the only parent if model = "s1" (or model = "s1pp").
p2ref	The reference counts for the second parent if model = "f1" (or model = "f1pp").
p2size	The total counts for the second parent if model = "f1" (or model = "f1pp").
snpname	A string. The name of the SNP under consideration. This is just returned in the input list for your reference.
bias_init	A vector of initial values for the bias parameter over the multiple runs of flexdog_full.
verbose	Should we output more (TRUE) or less (FALSE)?
outliers	A logical. Should we allow for the inclusion of outliers (TRUE) or not (FALSE). Only supported when model = "f1" or model = "s1". I wouldn't recommend it for any other model anyway.
prior_vec	The pre-specified genotype distribution. Only used if model = "custom" and must otherwise be NULL. If specified, then it should be a vector of length ploidy + 1 with non-negative elements that sum to 1.
...	Additional parameters to pass to flexdog_full.

**Details**

Possible values of the genotype distribution (values of model) are:

"norm" A distribution whose genotype frequencies are proportional to the density value of a normal with some mean and some standard deviation. Unlike the "bb" and "hw" options, this will allow for distributions both more and less dispersed than a binomial. This seems to be the most robust to violations in modeling assumptions, and so is the default. This prior class was developed in Gerard and Ferrão (2019).

"hw" A binomial distribution that results from assuming that the population is in Hardy-Weinberg equilibrium (HWE). This actually does pretty well even when there are minor to moderate deviations from HWE. Though it does not perform as well as the "norm" option when there are severe deviations from HWE.

- "bb" A beta-binomial distribution. This is an overdispersed version of "hw" and can be derived from a special case of the Balding-Nichols model.
- "ash" Any unimodal prior. This can sometimes be sensitive to violations in modeling assumptions, but tends to work better than the "flex" option. This prior class was developed in Gerard and Ferrão (2019).
- "s1" This prior assumes the individuals are all full-siblings resulting from one generation of selfing. I.e. there is only one parent. This model assumes a particular type of meiotic behavior: polysomic inheritance with bivalent, non-preferential pairing. Since this is a pretty strong and well-founded prior, we allow `outliers = TRUE` when `model = "s1"`.
- "s1pp" The same as "s1" but accounts for possible (and arbitrary levels of) preferential pairing during meiosis. The only supported values of `ploidy` right now for this option are 4 and 6.
- "f1" This prior assumes the individuals are all full-siblings resulting from one generation of a bi-parental cross. This model assumes a particular type of meiotic behavior: polysomic inheritance with bivalent, non-preferential pairing. Since this is a pretty strong and well-founded prior, we allow `outliers = TRUE` when `model = "f1"`.
- "f1pp" The same as "f1" but accounts for possible (and arbitrary levels of) preferential pairing during meiosis. This option is mostly untested for values of `ploidy` greater than 6.
- "flex" Generically any categorical distribution. Theoretically, this works well if you have a lot of individuals. In practice, it seems to be much less robust to violations in modeling assumptions.
- "uniform" A discrete uniform distribution. This should never be used in practice.
- "custom" A pre-specified prior distribution. You specify it using the `prior_vec` argument. You should almost never use this option in practice.

You might think a good default is `model = "uniform"` because it is somehow an "uninformative prior." But it is very informative and tends to work horribly in practice. The intuition is that it will estimate the allele bias and sequencing error rates so that the estimated genotypes are approximately uniform (since we are assuming that they are approximately uniform). This will usually result in unintuitive genotyping since most populations don't have a uniform genotype distribution. I include it as an option only for completeness. Please don't use it.

The value of `prop_mis` is a very intuitive measure for the quality of the SNP. `prop_mis` is the posterior proportion of individuals mis-genotyped. So if you want only SNPs that accurately genotype, say, 95% of the individuals, you could discard all SNPs with a `prop_mis` over 0.05.

The value of `maxpostprob` is a very intuitive measure for the quality of the genotype estimate of an individual. This is the posterior probability of correctly genotyping the individual when using `geno` (the posterior mode) as the genotype estimate. So if you want to correctly genotype, say, 95% of individuals, you could discard all individuals with a `maxpostprob` of under 0.95. However, if you are just going to impute missing genotypes later, you might consider not discarding any individuals as flexdog's genotype estimates will probably be more accurate than other more naive approaches, such as imputing using the grand mean.

In most datasets I've examined, allelic bias is a major issue. However, you may fit the model assuming no allelic bias by setting `update_bias = FALSE` and `bias_init = 1`.

Prior to using flexdog, during the read-mapping step, you could try to get rid of allelic bias by using WASP (<https://doi.org/10.1101/011221>). If you are successful in removing the allelic bias (because its only source was the read-mapping step), then setting `update_bias = FALSE` and `bias_init = 1` would be reasonable. You can visually inspect SNPs for bias by using `plot geno`.



flexdog, like most methods, is invariant to which allele you label as the "reference" and which you label as the "alternative". That is, if you set `refvec` with the number of alternative read-counts, then the resulting genotype estimates will be the estimated allele dosage of the alternative allele.

## Value

An object of class `flexdog`, which consists of a list with some or all of the following elements:

`bias` The estimated bias parameter.

`seq` The estimated sequencing error rate.

`od` The estimated overdispersion parameter.

`num_iter` The number of EM iterations ran. You should be wary if this equals `itermax`.

`llike` The maximum marginal log-likelihood.

`postmat` A matrix of posterior probabilities of each genotype for each individual. The rows index the individuals and the columns index the allele dosage.

`gene_dist` The estimated genotype distribution. The  $i$ th element is the proportion of individuals with genotype  $i-1$ . If `outliers = TRUE`, then this is conditional on the point not being an outlier.

`par` A list of the final estimates of the parameters of the genotype distribution. The elements included in `par` depends on the value of `model`:

`model = "norm"`:  $\mu$  is the normal mean and  $\sigma$  is the normal standard deviation (not variance).

`model = "hw"`:  $\alpha$  is the major allele frequency.

`model = "bb"`:  $\alpha$  is the major allele frequency and  $\tau$  is the overdispersion parameter (see the description of  $\rho$  in the Details of [betabinom](#)).

`model = "ash"`: `par` is an empty list.

`model = "s1"`: `pgeno` is the allele dosage of the parent and  $\alpha$  is the mixture proportion of the discrete uniform (included and fixed at a small value mostly for numerical stability reasons). See the description of `fs1_alpha` in [flexdog\\_full](#).

`model = "s1pp"`: `pgeno` is the allele dosage of the parent and `p1_pair_weights` contains a vector of mixing weights where element  $i$  is the mixing proportion for the segregation distribution in row  $i$  of `get_bivalent_probs(ploidy)$probmat[get_bivalent_probs(ploidy)$lvec == pgeno, , drop = FALSE]`.

`model = "f1"`: `p1geno` is the allele dosage of the first parent, `p2geno` is the allele dosage of the second parent, and  $\alpha$  is the mixture proportion of the discrete uniform (included and fixed at a small value mostly for numerical stability reasons). See the description of `fs1_alpha` in [flexdog\\_full](#).

`model = "f1pp"`: `p1geno` is the allele dosage of the first parent, `p2geno` is the allele dosage of the second parent, `p1_pair_weights` contains a vector of mixing weights where element  $i$  is the mixing proportion for the segregation distribution for parent 1 in row  $i$  of `get_bivalent_probs(ploidy)$probmat[get_bivalent_probs(ploidy)$lvec == p1geno, , drop = FALSE]`, and `p2_pair_weights` contains a vector of mixing weights where element  $i$  is the mixing proportion for the segregation distribution for parent 2 in row  $i$  of `get_bivalent_probs(ploidy)$probmat[get_bivalent_probs(ploidy)$lvec == p2geno, , drop = FALSE]`.

model = "flex": par is an empty list.  
 model = "uniform": par is an empty list.  
 model = "custom": par is an empty list.  
 geno The posterior mode genotype. These are your genotype estimates.  
 maxpostprob The maximum posterior probability. This is equivalent to the posterior probability of correctly genotyping each individual.  
 postmean The posterior mean genotype. In downstream association studies, you might want to consider using these estimates.  
 input\$refvec The value of refvec provided by the user.  
 input\$sizevec The value of sizevec provided by the user.  
 input\$ploidy The value of ploidy provided by the user.  
 input\$model The value of model provided by the user.  
 input\$p1ref The value of p1ref provided by the user.  
 input\$p1size The value of p1size provided by the user.  
 input\$p2ref The value of p2ref provided by the user.  
 input\$p2size The value of p2size provided by the user.  
 input\$snpname The value of snpname provided by the user.  
 prop\_mis The posterior proportion of individuals genotyped incorrectly.  
 out\_prop The estimated proportion of points that are outliers. Only available if outliers = TRUE.  
 prob\_out The ith element is the posterior probability that individual i is an outlier. Only available if outliers = TRUE.

### Author(s)

David Gerard

### References

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).
- Gerard, D. and Ferrão, L. F. V. (2019). Priors for Genotyping Polyploids. *bioRxiv*. doi: [10.1101/751784](https://doi.org/10.1101/751784).

### See Also

Run `browseVignettes(package = "updog")` in R for example usage. Other useful functions include:

[flexdog\\_full](#) For additional parameter options when running flexdog.

[rgeno](#) For simulating genotypes under the allowable prior models in flexdog.

[rflexdog](#) For simulating read-counts under the assumed likelihood model in flexdog.

[plot.flexdog](#) For plotting the output of flexdog.

[oracle\\_mis](#) For calculating the oracle genotyping error rates. This is useful for read-depth calculations *before* collecting data. After you have data, using the value of `prop_mis` is better.

[oracle\\_cor](#) For calculating the correlation between the true genotypes and an oracle estimator (useful for read-depth calculations *before* collecting data).

## Examples

```
## An S1 population where the first individual
## is the parent. Fit assuming outliers.
data("snpdat")
ploidy <- 6
refvec <- snpdat$counts[snpdat$snp == "SNP3"]
sizevec <- snpdat$size[snpdat$snp == "SNP3"]
fout <- flexdog(refvec = refvec[-1],
                sizevec = sizevec[-1],
                ploidy = ploidy,
                model = "s1",
                plref = refvec[1],
                plsize = sizevec[1],
                outliers = TRUE)

plot(fout)
```

```
## A natural population. We will assume a
## normal prior since there are so few
## individuals.
data("uitdewilligen")
ploidy <- 4
refvec <- uitdewilligen$refmat[, 1]
sizevec <- uitdewilligen$sizemat[, 1]
fout <- flexdog(refvec = refvec,
                sizevec = sizevec,
                ploidy = ploidy,
                model = "norm")

plot(fout)
```

---

flexdog\_full

*Flexible genotyping for polyploids from next-generation sequencing data.*


---

## Description

Genotype polyploid individuals from next generation sequencing (NGS) data while assuming the genotype distribution is one of several forms. `flexdog` does this while accounting for allele bias, overdispersion, sequencing error, and possibly outlying observations (if `model = "f1"` or `model = "s1"`). This function has more options than `flexdog` and is only meant for expert users. The method is described in detail in Gerard et. al. (2018) and Gerard and Ferrão (2019).

**Usage**

```
flexdog_full(refvec, sizevec, ploidy, model = c("norm", "hw", "bb",
"ash", "s1", "s1pp", "f1", "f1pp", "flex", "uniform", "custom"),
verbose = TRUE, mean_bias = 0, var_bias = 0.7^2, mean_seq = -4.7,
var_seq = 1, mean_od = -5.5, var_od = 0.5^2, seq = 0.005,
bias = 1, od = 0.001, update_bias = TRUE, update_seq = TRUE,
update_od = TRUE, mode = NULL, use_cvxr = FALSE, itermax = 200,
tol = 10^-4, fs1_alpha = 10^-3, ashpen = 10^-6, p1ref = NULL,
p1size = NULL, p2ref = NULL, p2size = NULL, snpname = NULL,
outliers = FALSE, prior_vec = NULL)
```

**Arguments**

refvec	A vector of counts of reads of the reference allele.
sizevec	A vector of total counts.
ploidy	The ploidy of the species. Assumed to be the same for each individual.
model	What form should the prior (genotype distribution) take? See Details for possible values.
verbose	Should we output more (TRUE) or less (FALSE)?
mean_bias	The prior mean of the log-bias.
var_bias	The prior variance of the log-bias.
mean_seq	The prior mean of the logit of the sequencing error rate.
var_seq	The prior variance of the logit of the sequencing error rate.
mean_od	The prior mean of the logit of the overdispersion parameter.
var_od	The prior variance of the logit of the overdispersion parameter.
seq	The starting value of the sequencing error rate.
bias	The starting value of the bias.
od	The starting value of the overdispersion parameter.
update_bias	A logical. Should we update bias (TRUE), or not (FALSE)?
update_seq	A logical. Should we update seq (TRUE), or not (FALSE)?
update_od	A logical. Should we update od (TRUE), or not (FALSE)?
mode	The mode if model = "ash". If not provided, flexdog will estimate the mode. This is the starting point of the allele frequency if model = "hw". This should be NULL for all other options of model.
use_cvxr	A logical. If model = "ash", then do you want to use the EM algorithm (FALSE) or a convex optimization program using the package CVXR (TRUE)? Only available if CVXR is installed. Setting use_cvxr to TRUE is generally slower than setting it to FALSE.
itermax	The maximum number of EM iterations to run for each mode (if model = "ash") or the total number of EM iterations to run (for any other value of model).
tol	The tolerance stopping criterion. The EM algorithm will stop if the difference in the log-likelihoods between two consecutive iterations is less than tol.

fs1_alpha	The value at which to fix the mixing proportion for the uniform component when model = "f1", model = "f1pp", model = "s1", or model = "s1pp". I would recommend some small value such as $10^{-3}$ .
ashpen	The penalty to put on the unimodal prior. Larger values shrink the unimodal prior towards the discrete uniform distribution.
p1ref	The reference counts for the first parent if model = "f1" (or model = "f1pp"), or for the only parent if model = "s1" (or model = "s1pp").
p1size	The total counts for the first parent if model = "f1" (or model = "f1pp"), or for the only parent if model = "s1" (or model = "s1pp").
p2ref	The reference counts for the second parent if model = "f1" (or model = "f1pp").
p2size	The total counts for the second parent if model = "f1" (or model = "f1pp").
snpname	A string. The name of the SNP under consideration. This is just returned in the input list for your reference.
outliers	A logical. Should we allow for the inclusion of outliers (TRUE) or not (FALSE). Only supported when model = "f1" or model = "s1". I wouldn't recommend it for any other model anyway.
prior_vec	The pre-specified genotype distribution. Only used if model = "custom" and must otherwise be NULL. If specified, then it should be a vector of length ploidy + 1 with non-negative elements that sum to 1.

## Details

Possible values of the genotype distribution (values of model) are:

- "norm" A distribution whose genotype frequencies are proportional to the density value of a normal with some mean and some standard deviation. Unlike the "bb" and "hw" options, this will allow for distributions both more and less dispersed than a binomial. This seems to be the most robust to violations in modeling assumptions, and so is the default. This prior class was developed in Gerard and Ferrão (2019).
- "hw" A binomial distribution that results from assuming that the population is in Hardy-Weinberg equilibrium (HWE). This actually does pretty well even when there are minor to moderate deviations from HWE. Though it does not perform as well as the "norm" option when there are severe deviations from HWE.
- "bb" A beta-binomial distribution. This is an overdispersed version of "hw" and can be derived from a special case of the Balding-Nichols model.
- "ash" Any unimodal prior. This can sometimes be sensitive to violations in modeling assumptions, but tends to work better than the "flex" option. This prior class was developed in Gerard and Ferrão (2019).
- "s1" This prior assumes the individuals are all full-siblings resulting from one generation of selfing. I.e. there is only one parent. This model assumes a particular type of meiotic behavior: polysomic inheritance with bivalent, non-preferential pairing. Since this is a pretty strong and well-founded prior, we allow outliers = TRUE when model = "s1".
- "s1pp" The same as "s1" but accounts for possible (and arbitrary levels of) preferential pairing during meiosis. The only supported values of ploidy right now for this option are 4 and 6.

"f1" This prior assumes the individuals are all full-siblings resulting from one generation of a bi-parental cross. This model assumes a particular type of meiotic behavior: polysomic inheritance with bivalent, non-preferential pairing. Since this is a pretty strong and well-founded prior, we allow outliers = TRUE when model = "f1".

"f1pp" The same as "f1" but accounts for possible (and arbitrary levels of) preferential pairing during meiosis. This option is mostly untested for values of ploidy greater than 6.

"flex" Generically any categorical distribution. Theoretically, this works well if you have a lot of individuals. In practice, it seems to be much less robust to violations in modeling assumptions.

"uniform" A discrete uniform distribution. This should never be used in practice.

"custom" A pre-specified prior distribution. You specify it using the prior\_vec argument. You should almost never use this option in practice.

You might think a good default is model = "uniform" because it is somehow an "uninformative prior." But it is very informative and tends to work horribly in practice. The intuition is that it will estimate the allele bias and sequencing error rates so that the estimated genotypes are approximately uniform (since we are assuming that they are approximately uniform). This will usually result in unintuitive genotyping since most populations don't have a uniform genotype distribution. I include it as an option only for completeness. Please don't use it.

The value of prop\_mis is a very intuitive measure for the quality of the SNP. prop\_mis is the posterior proportion of individuals mis-genotyped. So if you want only SNPs that accurately genotype, say, 95% of the individuals, you could discard all SNPs with a prop\_mis over 0.05.

The value of maxpostprob is a very intuitive measure for the quality of the genotype estimate of an individual. This is the posterior probability of correctly genotyping the individual when using geno (the posterior mode) as the genotype estimate. So if you want to correctly genotype, say, 95% of individuals, you could discard all individuals with a maxpostprob of under 0.95. However, if you are just going to impute missing genotypes later, you might consider not discarding any individuals as flexdog's genotype estimates will probably be more accurate than other more naive approaches, such as imputing using the grand mean.

In most datasets I've examined, allelic bias is a major issue. However, you may fit the model assuming no allelic bias by setting update\_bias = FALSE and bias\_init = 1.

Prior to using flexdog, during the read-mapping step, you could try to get rid of allelic bias by using WASP (<https://doi.org/10.1101/011221>). If you are successful in removing the allelic bias (because its only source was the read-mapping step), then setting update\_bias = FALSE and bias\_init = 1 would be reasonable. You can visually inspect SNPs for bias by using plot\_genom.

flexdog, like most methods, is invariant to which allele you label as the "reference" and which you label as the "alternative". That is, if you set refvec with the number of alternative read-counts, then the resulting genotype estimates will be the estimated allele dosage of the alternative allele.

## Value

An object of class flexdog, which consists of a list with some or all of the following elements:

bias The estimated bias parameter.

seq The estimated sequencing error rate.

od The estimated overdispersion parameter.

num\_iter The number of EM iterations ran. You should be wary if this equals itermax.

- llike The maximum marginal log-likelihood.
- postmat A matrix of posterior probabilities of each genotype for each individual. The rows index the individuals and the columns index the allele dosage.
- gene\_dist The estimated genotype distribution. The  $i$ th element is the proportion of individuals with genotype  $i-1$ . If `outliers = TRUE`, then this is conditional on the point not being an outlier.
- par A list of the final estimates of the parameters of the genotype distribution. The elements included in `par` depends on the value of `model`:
- `model = "norm"`: `mu` is the normal mean and `sigma` is the normal standard deviation (not variance).
  - `model = "hw"`: `alpha` is the major allele frequency.
  - `model = "bb"`: `alpha` is the major allele frequency and `tau` is the overdispersion parameter (see the description of `rho` in the Details of [betabinom](#)).
  - `model = "ash"`: `par` is an empty list.
  - `model = "s1"`: `pgeno` is the allele dosage of the parent and `alpha` is the mixture proportion of the discrete uniform (included and fixed at a small value mostly for numerical stability reasons). See the description of `fs1_alpha` in [flexdog\\_full](#).
  - `model = "s1pp"`: `pgeno` is the allele dosage of the parent and `p1_pair_weights` contains a vector of mixing weights where element  $i$  is the mixing proportion for the segregation distribution in row  $i$  of `get_bivalent_probs(ploidy)$probsmat[get_bivalent_probs(ploidy)$lvec == pgeno, , drop = FALSE]`.
  - `model = "f1"`: `p1geno` is the allele dosage of the first parent, `p2geno` is the allele dosage of the second parent, and `alpha` is the mixture proportion of the discrete uniform (included and fixed at a small value mostly for numerical stability reasons). See the description of `fs1_alpha` in [flexdog\\_full](#).
  - `model = "f1pp"`: `p1geno` is the allele dosage of the first parent, `p2geno` is the allele dosage of the second parent, `p1_pair_weights` contains a vector of mixing weights where element  $i$  is the mixing proportion for the segregation distribution for parent 1 in row  $i$  of `get_bivalent_probs(ploidy)$probsmat[get_bivalent_probs(ploidy)$lvec == p1geno, , drop = FALSE]`, and `p2_pair_weights` contains a vector of mixing weights where element  $i$  is the mixing proportion for the segregation distribution for parent 2 in row  $i$  of `get_bivalent_probs(ploidy)$probsmat[get_bivalent_probs(ploidy)$lvec == p2geno, , drop = FALSE]`.
  - `model = "flex"`: `par` is an empty list.
  - `model = "uniform"`: `par` is an empty list.
  - `model = "custom"`: `par` is an empty list.
- geno The posterior mode genotype. These are your genotype estimates.
- maxpostprob The maximum posterior probability. This is equivalent to the posterior probability of correctly genotyping each individual.
- postmean The posterior mean genotype. In downstream association studies, you might want to consider using these estimates.
- input\$refvec The value of `refvec` provided by the user.
- input\$sizevec The value of `sizevec` provided by the user.
- input\$ploidy The value of `ploidy` provided by the user.

`input$model` The value of `model` provided by the user.  
`input$p1ref` The value of `p1ref` provided by the user.  
`input$p1size` The value of `p1size` provided by the user.  
`input$p2ref` The value of `p2ref` provided by the user.  
`input$p2size` The value of `p2size` provided by the user.  
`input$snpname` The value of `snpname` provided by the user.  
`prop_mis` The posterior proportion of individuals genotyped incorrectly.  
`out_prop` The estimated proportion of points that are outliers. Only available if `outliers = TRUE`.  
`prob_out` The *i*th element is the posterior probability that individual *i* is an outlier. Only available if `outliers = TRUE`.

### Author(s)

David Gerard

### References

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).
- Gerard, D. and Ferrão, L. F. V. (2019). Priors for Genotyping Polyploids. *bioRxiv*. doi: [10.1101/751784](https://doi.org/10.1101/751784).

### See Also

Run `browseVignettes(package = "updog")` in R for example usage. Other useful functions include:

[flexdog](#) For a more user-friendly version of `flexdog_full`.

[rgeno](#) For simulating genotypes under the allowable prior models in `flexdog`.

[rflexdog](#) For simulating read-counts under the assumed likelihood model in `flexdog`.

[plot.flexdog](#) For plotting the output of `flexdog`.

[oracle\\_mis](#) For calculating the oracle genotyping error rates. This is useful for read-depth calculations *before* collecting data. After you have data, using the value of `prop_mis` is better.

[oracle\\_cor](#) For calculating the correlation between the true genotypes and an oracle estimator (useful for read-depth calculations *before* collecting data).

### Examples

```

## A natural population. We will assume a
## normal prior since there are so few
## individuals.
data("uitdewilligen")
ploidy <- 4
refvec <- uitdewilligen$refmat[, 1]
sizevec <- uitdewilligen$sizemat[, 1]
fout <- flexdog_full(refvec = refvec,

```



```

                                sizevec = sizevec,
                                ploidy  = ploidy,
                                model   = "norm")
plot(fout)

```

---

flexdog\_obj                      *Log-likelihood that [flexdog](#) maximizes.*

---

### Description

Log-likelihood that [flexdog](#) maximizes.

### Usage

```
flexdog_obj(prob_k_vec, ref_vec, size_vec, ploidy, seq, bias, od, mean_bias,
            var_bias, mean_seq, var_seq, mean_od, var_od)
```

### Arguments

prob_k_vec	The kth element is the prior probability of genotype k (when starting to count from 0).
ref_vec	A vector of counts of reads of the reference allele.
size_vec	A vector of total counts.
ploidy	The ploidy of the species. Assumed to be the same for each individual.
seq	The starting value of the sequencing error rate.
bias	The starting value of the bias.
od	The starting value of the overdispersion parameter.
mean_bias	The prior mean of the log-bias.
var_bias	The prior variance of the log-bias.
mean_seq	The prior mean of the logit of the sequencing error rate.
var_seq	The prior variance of the logit of the sequencing error rate.
mean_od	The prior mean of the logit of the overdispersion parameter.
var_od	The prior variance of the logit of the overdispersion parameter.

### Value

The objective (marginal log-likelihood) used in [flexdog\\_full](#).

### Author(s)

David Gerard

---

flexdog\_obj\_out      *Log-likelihood that flexdog maximizes when outliers are present.*

---

### Description

Log-likelihood that flexdog maximizes when outliers are present.

### Usage

```
flexdog_obj_out(prob_k_vec, out_prop, ref_vec, size_vec, ploidy, seq, bias,
               od, mean_bias, var_bias, mean_seq, var_seq, mean_od, var_od)
```

### Arguments

prob_k_vec	The kth element is the prior probability of genotype k (when starting to count from 0).
out_prop	The probability of being an outlier.
ref_vec	A vector of counts of reads of the reference allele.
size_vec	A vector of total counts.
ploidy	The ploidy of the species. Assumed to be the same for each individual.
seq	The starting value of the sequencing error rate.
bias	The starting value of the bias.
od	The starting value of the overdispersion parameter.
mean_bias	The prior mean of the log-bias.
var_bias	The prior variance of the log-bias.
mean_seq	The prior mean of the logit of the sequencing error rate.
var_seq	The prior variance of the logit of the sequencing error rate.
mean_od	The prior mean of the logit of the overdispersion parameter.
var_od	The prior variance of the logit of the overdispersion parameter.

### Value

A double. The flexdog objective when outliers = TRUE.

### Author(s)

David Gerard

### See Also

[flexdog\\_obj](#) for the objective function without outliers.

---

flex\_update\_pivec      *Update the distribution of genotypes from various models.*

---

### Description

Update the distribution of genotypes from various models.

### Usage

```
flex_update_pivec(weight_vec, model = c("hw", "bb", "norm", "ash", "f1",
    "s1", "f1pp", "s1pp", "f1ppdr", "s1ppdr", "flex", "uniform", "custom"),
    control)
```

### Arguments

weight_vec	colSums(wik_mat) from <a href="#">flexdog</a> . This is the sum of current posterior probabilities of each individual having genotype k.
model	What model are we assuming? See the description in <a href="#">flexdog</a> for details.
control	A list of anything else needed to be passed. E.g. if model = "ash", then inner_weights needs to be passed through control (see <a href="#">get_inner_weights</a> for how to get this matrix).

### Value

A list with the following elements

pivec The estimate of the genotype distribution.

par A list of estimated parameters. An empty list if the model does not contain any parameters other than pivec.

### Author(s)

David Gerard

---

get\_bivalent\_probs      *Returns segregation probabilities, pairing representation and number of ref alleles given the ploidy.*

---

### Description

Returns segregation probabilities, pairing representation and number of ref alleles given the ploidy.

### Usage

```
get_bivalent_probs(ploidy)
```

**Arguments**

`ploidy` The ploidy of the individual. Should be even and greater than 0.

**Value**

A list of three elements

`propmat` The rows index the pairing configuration and the columns index the number of reference alleles segregating. The elements are the probability of segregating the given number of reference alleles in a given category.

`pmat` The pairing representation of the configuration.

`lvec` The number of reference alleles an individual has given their pairing configuration in `pmat`.

**Author(s)**

David Gerard

**Examples**

```
get_bivalent_probs(4)
```

---

`get_bivalent_probs_dr` *Double reduction version of [get\\_bivalent\\_probs](#).*

---

**Description**

Double reduction version of [get\\_bivalent\\_probs](#).

**Usage**

```
get_bivalent_probs_dr(ploidy)
```

**Arguments**

`ploidy` The ploidy of the individual. Should be even and greater than 0.

**Value**

A list. The same elements as in [get\\_bivalent\\_probs](#), augmented to include more scenarios, but with the additional element `penvec`. This is a logical vector that is TRUE if the corresponding rows of `propmat` and `pmat` and elements of `lvec` would be included in the non-double-reduction-model and is FALSE otherwise.

**Author(s)**

David Gerard

**See Also**

[get\\_bivalent\\_probs](#).

**Examples**

```
get_bivalent_probs(4)
```

---

```
get_conv_inner_weights
```

*Get the inner weights used for the em update in [update\\_pp\\_f1](#) when there are more than two bivalent components for one of the parents.*

---

**Description**

Get the inner weights used for the em update in [update\\_pp\\_f1](#) when there are more than two bivalent components for one of the parents.

**Usage**

```
get_conv_inner_weights(psegprob, psegmat)
```

**Arguments**

<code>psegprob</code>	One of the parents segregation probability vector.
<code>psegmat</code>	The other parent's segregation matrix.

**Value**

A matrix. The columns index the K components (aka individuals in the context of the local problem) and the rows index the bivalent components.

**Author(s)**

David Gerard

**See Also**

[update\\_pp\\_f1](#) for where this is used. [uni\\_em\\_const](#) for where the weights are used (equivalent to `lmat` there).

---

get_dimname	<i>Returns a vector character strings that are all of the possible combinations of the reference allele and the non-reference allele.</i>
-------------	---

---

**Description**

Returns a vector character strings that are all of the possible combinations of the reference allele and the non-reference allele.

**Usage**

```
get_dimname(ploidy)
```

**Arguments**

ploidy            The ploidy of the species.

**Value**

For example, if ploidy = 3 then this will return `c("aaa", "Aaa", "AAa", "AAA")`

**Author(s)**

David Gerard

---

get_hyper_weights	<i>Return mixture weights needed to obtain a hypergeometric distribution.</i>
-------------------	---

---

**Description**

Obtains the mixing weights for the mixing distributions of `get_bivalent_probs` to return a hypergeometric distribution where ploidy is the population size, e11 is the number of success states in the population, and ploidy / 2 is the number of draws. If these are the mixing weights in the population, then there is no preferential pairing.

**Usage**

```
get_hyper_weights(ploidy, e11)
```

**Arguments**

ploidy            The ploidy of the individual.  
e11                The number of reference alleles in the individual.

**Value**

A list with the following two elements:

`pmat` Each row is a category and the columns index either aa, Aa, or AA.

`weightvec` The mixing weights for each row of `pmat`.

**Author(s)**

David Gerard

**Examples**

```
get_hyper_weights(4, 2)
```

---

<code>get_inner_weights</code>	<i>Compute inner weights for updating the mixing proportions when using ash model.</i>
--------------------------------	--

---

**Description**

The (i,k)th element is  $1(k \in F(a, i))/|F(a, i)|$ .

**Usage**

```
get_inner_weights(ploidy, mode)
```

**Arguments**

`ploidy` The ploidy of the species. Assumed to be the same for each individual.

`mode` The mode if `model = "ash"`. If not provided, `flexdog` will estimate the mode. This is the starting point of the allele frequency if `model = "hw"`. This should be NULL for all other options of `model`.

**Value**

A matrix of numerics. The weights used for the weighted EM algorithm in [flexdog\\_full](#).

**Author(s)**

David Gerard

---

get_probk_vec	<i>Obtain the genotype distribution given the distribution of discrete uniforms.</i>
---------------	--

---

**Description**

Obtain the genotype distribution given the distribution of discrete uniforms.

**Usage**

```
get_probk_vec(pivec, model, mode)
```

**Arguments**

pivec	The mixing probability of the i'th discrete uniform distribution.
model	What form should the prior (genotype distribution) take? See Details for possible values.
mode	The mode if model = "ash". If not provided, flexdog will estimate the mode. This is the starting point of the allele frequency if model = "hw". This should be NULL for all other options of model.

**Value**

A vector of numerics. Element k is the probability of genotype k.

**Author(s)**

David Gerard

**See Also**

[flexdog](#) where this is used.

---

get_q_array	<i>Return the probabilities of an offspring's genotype given its parental genotypes for all possible combinations of parental and offspring genotypes. This is for species with polysomal inheritance and bivalent, non-preferential pairing.</i>
-------------	---

---

**Description**

Return the probabilities of an offspring's genotype given its parental genotypes for all possible combinations of parental and offspring genotypes. This is for species with polysomal inheritance and bivalent, non-preferential pairing.



**Usage**

```
get_q_array(ploidy)
```

**Arguments**

ploidy            A positive integer. The ploidy of the species.

**Value**

An three-way array of proportions. The (i, j, k)th element is the probability of an offspring having k - 1 reference alleles given that parent 1 has i - 1 reference alleles and parent 2 has j - 1 reference alleles. Each dimension of the array is ploidy + 1. In the dimension names, "A" stands for the reference allele and "a" stands for the alternative allele.

**Author(s)**

David Gerard

**Examples**

```
qarray <- get_q_array(6)
apply(qarray, c(1, 2), sum) ## should all be 1's.
```

---

get\_uni\_rep            *Get the representation of a discrete unimodal probability distribution.*

---

**Description**

NB: In `get_uni_rep`, we count the mode starting at 1. In [get\\_probk\\_vec](#), we count the mode starting at 0.

**Usage**

```
get_uni_rep(probvec)
```

**Arguments**

probvec            A probability vector. It assumes the probabilities are ordered according to the ordering of the discrete set.

**Value**

A list with the following elements

pivec The mixing weights for the unimodal representation.

mode The central value of the unimodal distribution.

**Author(s)**

David Gerard

**See Also**

[get\\_probk\\_vec](#) with option `model = "ash"` for the inverse of this function.

---

get\_wik\_mat

*E-step in [flexdog](#).*

---

**Description**

E-step in [flexdog](#).

**Usage**

```
get_wik_mat(prob_k_vec, ref_vec, size_vec, ploidy, seq, bias, od)
```

**Arguments**

<code>prob_k_vec</code>	The vector of current prior probabilities of each genotype.
<code>ref_vec</code>	A vector of counts of reads of the reference allele.
<code>size_vec</code>	A vector of total counts.
<code>ploidy</code>	The ploidy of the species. Assumed to be the same for each individual.
<code>seq</code>	The starting value of the sequencing error rate.
<code>bias</code>	The starting value of the bias.
<code>od</code>	The starting value of the overdispersion parameter.

**Value**

A matrix of numerics. The rows index the individuals and the columns index the genotype. These weights are used in the EM algorithm (and is indeed the E-step) in [flexdog\\_full](#).

**Author(s)**

David Gerard

**See Also**

[flexdog](#) for the full EM algorithm.

---

get\_wik\_mat\_out      *E-step in flexdog where we now allow an outlier distribution.*

---

### Description

E-step in [flexdog](#) where we now allow an outlier distribution.

### Usage

```
get_wik_mat_out(probkc_vec, out_prop, refvec, sizevec, ploidy, seq, bias,
od)
```

### Arguments

probkc_vec	The vector of current prior probabilities of each genotype.
out_prop	The probability of being an outlier.
refvec	A vector of counts of reads of the reference allele.
sizevec	A vector of total counts.
ploidy	The ploidy of the species. Assumed to be the same for each individual.
seq	The starting value of the sequencing error rate.
bias	The starting value of the bias.
od	The starting value of the overdispersion parameter.

### Value

Same as [get\\_wik\\_mat](#) but the last column is for the outlier class.

### Author(s)

David Gerard

### See Also

[flexdog](#) for the full EM algorithm. [get\\_wik\\_mat](#) for the equivalent function without outliers. [doutdist](#) for the outlier density function.

---

grad\_for\_eps                      *Gradient for obj\_for\_eps.*

---

### Description

Gradient for [obj\\_for\\_eps](#).

### Usage

```
grad_for_eps(parvec, refvec, sizevec, ploidy, mean_bias, var_bias,
             mean_seq, var_seq, mean_od, var_od, wmat, update_bias = TRUE,
             update_seq = TRUE, update_od = TRUE)
```

### Arguments

parvec	A vector of length three. The first element is the sequencing error rate, the second element is the allele bias, and the third element is the overdispersion parameter.
refvec	A vector. The <i>i</i> th element is the reference count for the <i>i</i> th individual in the SNP.
sizevec	A vector. the <i>i</i> th element is the size count for the <i>i</i> th individual in the SNP/
ploidy	The ploidy of the species.
mean_bias	The prior mean of the log-bias.
var_bias	The prior variance of the log-bias
mean_seq	The prior mean of the logit sequencing error rate.
var_seq	The prior variance of the logit sequencing error rate.
mean_od	The prior mean of the logit of the overdispersion parameter
var_od	The prior variance of the logit of the overdispersion parameter.
wmat	The matrix of (variational) posterior probabilities for each dosage. The rows index the individuals and the columns index the dosage levels.
update_bias	A logical. This is not used in <a href="#">obj_for_eps</a> , but sets the second element to 0.0 in <a href="#">grad_for_eps</a> .
update_seq	A logical. This is not used in <a href="#">obj_for_eps</a> , but sets the first element to 0.0 in <a href="#">grad_for_eps</a> .
update_od	A logical. This is not used in <a href="#">obj_for_eps</a> , but sets the third element to 0.0 in <a href="#">grad_for_eps</a> .

### Value

A double.

### Author(s)

David Gerard

---

grad\_for\_mu\_sigma2      *Gradient for [obj\\_for\\_mu\\_sigma2](#) with respect for mu and sigma2.*

---

### Description

Gradient for [obj\\_for\\_mu\\_sigma2](#) with respect for mu and sigma2.

### Usage

```
grad_for_mu_sigma2(mu, sigma2, phifk_mat, cor_inv, log_bb_dense)
```

### Arguments

mu	A vector, the ith element is the variational posterior mean of individual i for the SNP.
sigma2	A vector, the ith element is the variational posterior variance of individual i for the SNP.
phifk_mat	A matrix that contains the standard normal quantile of the beta-binomial cdf at dosage k for individual i. The rows index the individuals and the columns index the dosages.
cor_inv	The inverse of the underlying correlation matrix.
log_bb_dense	A matrix of log-densities of the beta binomial. The rows index the individuals and the columns index the allele dosage. Allele dosage goes from -1 to ploidy, so there are ploidy + 2 elements.

### Value

A vector of length  $2 * n_{ind}$  of numerics. The first element n elements are the partial derivatives with respect to mu and the second n elements are the partial derivatives with respect to sigma2 in [obj\\_for\\_mu\\_sigma2](#).

### Author(s)

David Gerard

---

grad\_for\_mu\_sigma2\_wrapper  
*Gradient for [obj\\_for\\_mu\\_sigma2\\_wrapper](#) with respect for muSigma2 and a wrapper for [grad\\_for\\_mu\\_sigma2](#)*

---

### Description

Gradient for [obj\\_for\\_mu\\_sigma2\\_wrapper](#) with respect for muSigma2 and a wrapper for [grad\\_for\\_mu\\_sigma2](#)

**Usage**

```
grad_for_mu_sigma2_wrapper(muSigma2, phifk_mat, cor_inv, log_bb_dense)
```

**Arguments**

muSigma2	A vector. The first half are mu and the second half are sigma2.
phifk_mat	A matrix that contains the standard normal quantile of the beta-binomial cdf at dosage k for individual i. The rows index the individuals and the columns index the dosages.
cor_inv	The inverse of the underlying correlation matrix.
log_bb_dense	A matrix of log-densities of the beta binomial. The rows index the individuals and the columns index the allele dosage. Allele dosage goes from -1 to ploidy, so there are ploidy + 2 elements.

**Value**

A vector of length  $2 * n_{ind}$  of numerics. The first element n elements are the partial derivatives with respect to mu and the second n elements are the partial derivatives with respect to sigma2 in [obj\\_for\\_mu\\_sigma2](#).

**Author(s)**

David Gerard

---

grad\_for\_weighted\_lbb *Gradient for [obj\\_for\\_weighted\\_lbb](#).*

---

**Description**

Gradient for [obj\\_for\\_weighted\\_lbb](#).

**Usage**

```
grad_for_weighted_lbb(parvec, ploidy, weight_vec)
```

**Arguments**

parvec	A vector of length 2. The first term is the current mean of the underlying beta. The second term is the current overdispersion parameter.
ploidy	The ploidy of the species.
weight_vec	A vector of length ploidy + 1 that contains the weights for each component beta-binomial.

**Value**

A vector of length 2. The first component is the gradient of the mean of the underlying beta. The second component is the gradient of the overdispersion parameter of the underlying beta.

**Author(s)**

David Gerard

---

`grad_for_weighted_lnorm`*Gradient for [obj\\_for\\_weighted\\_lnorm](#).*

---

**Description**Gradient for [obj\\_for\\_weighted\\_lnorm](#).**Usage**`grad_for_weighted_lnorm(parvec, ploidy, weight_vec)`**Arguments**

<code>parvec</code>	A vector of length 2. The first term is the current mean of the underlying normal. The second term is the current standard deviation (not variance) of the normal.
<code>ploidy</code>	The ploidy of the species.
<code>weight_vec</code>	A vector of length <code>ploidy + 1</code> that contains the weights for each component beta-binomial.

**Value**

A vector of length 2. The first term is the derivative with respect to the mean, the second term is the derivative with respect to the standard deviation (not variance).

**Author(s)**

David Gerard

---

`initialize_pivec`*Initialize pivec for [flexdog](#) EM algorithm.*

---

**Description**

The key idea here is choosing the pi's so that the two modes have equal probability.

**Usage**`initialize_pivec(ploidy, mode, model = c("hw", "bb", "norm", "ash", "f1", "s1", "f1pp", "s1pp", "f1ppdr", "s1ppdr", "flex", "uniform", "custom"))`

**Arguments**

ploidy	The ploidy of the species. Assumed to be the same for each individual.
mode	The mode if model = "ash". If not provided, flexdog will estimate the mode. This is the starting point of the allele frequency if model = "hw". This should be NULL for all other options of model.
model	What form should the prior (genotype distribution) take? See Details for possible values.

**Value**

A vector of numerics. The initial value of pivec used in `flexdog_full`.

**Author(s)**

David Gerard

**See Also**

`flexdog` for where this is used.

---

is.flexdog	<i>Tests if an argument is a flexdog object.</i>
------------	--

---

**Description**

Tests if an argument is a flexdog object.

**Usage**

```
is.flexdog(x)
```

**Arguments**

x	Anything.
---	-----------

**Value**

A logical. TRUE if x is a flexdog object, and FALSE otherwise.

**Author(s)**

David Gerard

**Examples**

```
is.flexdog("anything")  
# FALSE
```



---

is.mupdog                      *Tests if its argument is a mupdog object.*

---

**Description**

Tests if its argument is a mupdog object.

**Usage**

```
is.mupdog(x)
```

**Arguments**

x                      Anything.

**Value**

A logical. TRUE if x is a mupdog object, and FALSE otherwise.

**Author(s)**

David Gerard

**Examples**

```
is.mupdog("anything")  
# FALSE
```

---

logit                      *The logit function.*

---

**Description**

The logit function.

**Usage**

```
logit(x)
```

**Arguments**

x                      A double between 0 and 1.

**Value**

The logit of x.

**Author(s)**

David Gerard

---

log_sum_exp	<i>Log-sum-exponential trick.</i>
-------------	-----------------------------------

---

**Description**

Log-sum-exponential trick.

**Usage**

log\_sum\_exp(x)

**Arguments**

x                    A vector to log-sum-exp.

**Value**

The log of the sum of the exponential of the elements in x.

**Author(s)**

David Gerard

---

log_sum_exp_2	<i>Log-sum-exponential trick using just two doubles.</i>
---------------	--

---

**Description**

Log-sum-exponential trick using just two doubles.

**Usage**

log\_sum\_exp\_2(x, y)

**Arguments**x                    A double.  
y                    Another double.**Value**

The log of the sum of the exponential of x and y.

**Author(s)**

David Gerard

---

mupdog *Multi-SNP updog.*

---

### Description

A method to genotype autopolyploids using GBS or RAD-seq like data by accounting for correlations in the genotype distribution between the individuals.

### Usage

```
mupdog(refmat, sizemat, ploidy, verbose = TRUE, mean_bias = 0,
       var_bias = 1, mean_seq = -4.7, var_seq = 1, seq = NULL,
       bias = NULL, od = NULL, allele_freq = NULL, inbreeding = NULL,
       cor_mat = NULL, postmean = NULL, postvar = NULL,
       update_cor = TRUE, update_inbreeding = TRUE,
       update_allele_freq = TRUE, num_core = 1, update_method = c("Brent",
       "L-BFGS-B"), control = list())
```

### Arguments

refmat	A matrix of reference counts. The rows index the individuals and the columns index the SNPs.
sizemat	A matrix of total counts. The rows index the individuals and the columns index the SNPs. Should have the same dimensions as refmat.
ploidy	The ploidy of the species. To estimate the ploidy, re-run mupdog at various ploidy levels and choose the one with the largest ELBO. This assumes that the ploidy is the same for all individuals in the sample.
verbose	Should we print a lot of output TRUE or not FALSE?
mean_bias	The prior mean of the log-bias. Defaults to 0 (no bias).
var_bias	The prior variance on the log-bias. Defaults to 1. This roughly corresponds to likely bias values between 0.14 and 7.4. This is a far wider interval than what we observe in practice, thus making this prior rather uninformative. We usually observe bias values somewhere between 0.5 and 2.
mean_seq	The prior mean of the logit-sequencing-error-rate. Defaults to -4.7. This corresponds to a sequencing error rate of 0.009.
var_seq	The prior variance of the logit-sequencing-error-rate. Defaults to 1. This corresponds to likely values of 0.001 and 0.06. This upper bound is larger than what we would expect given the current state of next-gen-sequencing technology.
seq	A vector of initial sequencing errors. Should be the same length as the number of columns of refmat (number of SNPs). Must be between 0 and 1.
bias	A vector of initial bias parameters. Should be the same length as the number of columns of refmat (number of SNPs). Must be greater than 0.
od	A vector of initial overdispersion parameters. Should be the same length as the number of columns of refmat (number of SNPs). Must be between 0 and 1.

allele_freq	A vector of initial allele frequencies. Should be the same length as the number of columns of refmat (number of SNPs). Must be between 0 and 1.
inbreeding	A vector of initial individual-specific inbreeding coefficients. Should be the same length as the number of rows of refmat (number of individuals). Must be between 0 and 1.
cor_mat	Initial correlation matrix. Should have the same number of columns/rows as the number of individuals.
postmean	Initial variational posterior means. Should have the same dimensions as refmat.
postvar	Initial posterior variances. Should have the same dimensions as refmat.
update_cor	A logical. Should we update the underlying correlation matrix TRUE or not FALSE. It might be unwise to set this to TRUE if you have more individuals than SNPs.
update_inbreeding	A logical. Should we update the inbreeding coefficients TRUE or not FALSE?
update_allele_freq	A logical. Should we update the allele frequencies TRUE or not FALSE?
num_core	The number of cores to use if you want to run the optimization steps in parallel. If num_core = 1, then the optimization step will not be run in parallel.
update_method	What generic optimizer should we use to update allele_freq and inbreeding? Options are either "Brent" or "L-BFGS-B". See <a href="#">optim</a> for details on these optimizers.
control	A list of control parameters (itermax, obj_tol).

## Details

Blischak et al (2017) developed a genotyping approach for autopolyploids that assumes a Balding-Nichols generative model (Balding and Nichols, 1997) on the genotypes. Using a different generative model, Gerard et al (2018) accounted for common issues in sequencing data ignored by previous researchers. Mupdog unites and extends these two approaches:

- Unite: We account for locus-specific allele-bias, locus-specific sequencing error, and locus-specific overdispersion while marginally assuming a Balding-Nichols generative model on the genotypes.
- Extend: We account for underlying correlations between the individuals using a Gaussian copula model.

Mupdog uses a variational Bayes approach to estimate all parameters of interest and the posterior probabilities of the genotypes for each individual at each locus.

## Value

A list with some or all of the following elements:

map\_dosage A matrix of numerics containing the variational maximum-a-posterior (MAP) genotypes (allele dosages) for each individual at each SNP. Element (i, j) is the MAP genotype for individual i at SNP j.

**maxpostprob** A matrix of numerics containing the variational maximum posterior probabilities for each individual at each SNP. The (i, j) element is the variational posterior probability that individual i is genotyped correctly at SNP j.  
**postprob** A three-way array of numerics. Element (i, j, k) is the variational posterior probability that individual i has genotype k-1 at SNP j.  
**seq** A vector of numerics. Element j is the estimated sequencing error rate for SNP j.  
**bias** A vector of numerics. Element j is the estimated allelic bias for SNP j.  
**od** A vector of numerics. Element j is the estimated overdispersion parameter for SNP j.  
**allele\_freq** A vector of numerics. Element j is the estimated allele-frequency for SNP j.  
**inbreeding** A vector of numerics. Element i is the estimated inbreeding coefficient for individual i.  
**cor\_mat** A symmetric matrix of numerics. Element (i, j) is the estimated `_latent_` correlation between individual i and individual j.  
**postmean** A matrix of numerics. Element (i, j) is the variational posterior mean for individual i at SNP j.  
**postvar** A matrix of numerics. Element (i, j) is the variational posterior variance for individual i at SNP j.  
**input\$refmat** A matrix of numerics. The inputted `refmat`.  
**input\$sizemat** A matrix of numerics. The inputted `sizemat`.  
**input\$ploidy** The inputted `ploidy`.  
**obj** The maximized variational objective.

### Author(s)

David Gerard

### References

- David J Balding and Richard A Nichols. **Significant genetic correlations among caucasians at forensic DNA loci**. *Heredity*, 78(6):583–589, 1997. doi: 10.1038/sj.hdy.6881750.
- Paul D Blischak, Laura S Kubatko, and Andrea D Wolfe. **SNP genotyping and parameter estimation in polyploids using low-coverage sequencing data**. *Bioinformatics*, page btx587, 2017. doi: 10.1093/bioinformatics/btx587.
- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

### Examples

```

data(uitdewilligen)
mout <- mupdog(refmat = uitdewilligen$refmat,
               sizemat = uitdewilligen$sizemat,
               ploidy = uitdewilligen$ploidy,
               verbose = FALSE,
               control = list(obj_tol = 10^-4))

```

```
## Summaries of output
plot(mout, 4)
hist(mout$bias)
hist(mout$seq)
hist(mout$od)
hist(mout$inbreeding)
hist(mout$allele_freq)

## mupdog can correctly estimate ploidy to be 4
mout2 <- mupdog(refmat = uitdewilligen$refmat,
               sizemat = uitdewilligen$sizemat,
               ploidy = 2,
               verbose = FALSE,
               control = list(obj_tol = 10^-4))

mout6 <- mupdog(refmat = uitdewilligen$refmat,
               sizemat = uitdewilligen$sizemat,
               ploidy = 6,
               verbose = FALSE,
               control = list(obj_tol = 10^-4))

mout8 <- mupdog(refmat = uitdewilligen$refmat,
               sizemat = uitdewilligen$sizemat,
               ploidy = 8,
               verbose = FALSE,
               control = list(obj_tol = 10^-4))

y <- c(mout2$obj, mout$obj, mout6$obj, mout8$obj)
x <- seq(2, 8, by = 2)
plot(x, y, type = "l", xlab = "ploidy", ylab = "objective")
```

---

mupout

*A mupdog fit of the [uitdewilligen](#) data.*

---

### Description

A mupdog fit of the [uitdewilligen](#) data.

### Usage

mupout

### Format

An object of class [mupdog](#).

**Value**

See the Format Section.

**Source**

The raw data that this was fit to can be found in [uitdewilligen](#).

**See Also**

[uitdewilligen](#) The raw data.

[plot.mupdog](#) A method to plot a [mupdog](#) object.

[summary.mupdog](#) Calculate some summaries of a [mupdog](#) object.

[mupdog](#) Function used to create this [mupdog](#) object.

---

 obj\_for\_alpha

*Objective function when updating alpha*


---

**Description**

Objective function when updating alpha

**Usage**

```
obj_for_alpha(mu, sigma2, alpha, rho, log_bb_dense, ploidy)
```

**Arguments**

mu	A vector. The ith element is individual i's variational posterior mean at the SNP.
sigma2	A vector. The ith element is individual i's variational posterior variance at the SNP.
alpha	The SNP's allele frequency.
rho	A vector. The ith element is individuals i's inbreeding coefficient.
log_bb_dense	A matrix of log-densities of the beta binomial. The rows index the individuals and the columns index the allele dosage.
ploidy	The ploidy of the species.

**Value**

A double. The objective when updating alpha in [mupdog](#).

**Author(s)**

David Gerard

---

obj_for_eps	<i>Objective function for updating sequencing error rate, bias, and overdispersion parameters.</i>
-------------	--

---

### Description

Objective function for updating sequencing error rate, bias, and overdispersion parameters.

### Usage

```
obj_for_eps(parvec, refvec, sizevec, ploidy, mean_bias, var_bias, mean_seq,
  var_seq, mean_od, var_od, wmat, update_bias = TRUE,
  update_seq = TRUE, update_od = TRUE)
```

### Arguments

parvec	A vector of length three. The first element is the sequencing error rate, the second element is the allele bias, and the third element is the overdispersion parameter.
refvec	A vector. The <i>i</i> th element is the reference count for the <i>i</i> th individual in the SNP.
sizevec	A vector. the <i>i</i> th element is the size count for the <i>i</i> th individual in the SNP/
ploidy	The ploidy of the species.
mean_bias	The prior mean of the log-bias.
var_bias	The prior variance of the log-bias
mean_seq	The prior mean of the logit sequencing error rate.
var_seq	The prior variance of the logit sequencing error rate.
mean_od	The prior mean of the logit of the overdispersion parameter
var_od	The prior variance of the logit of the overdispersion parameter.
wmat	The matrix of (variational) posterior probabilities for each dosage. The rows index the individuals and the columns index the dosage levels.
update_bias	A logical. This is not used in <code>obj_for_eps</code> , but sets the second element to 0.0 in <code>grad_for_eps</code> .
update_seq	A logical. This is not used in <code>obj_for_eps</code> , but sets the first element to 0.0 in <code>grad_for_eps</code> .
update_od	A logical. This is not used in <code>obj_for_eps</code> , but sets the third element to 0.0 in <code>grad_for_eps</code> .

### Value

A double. The objective when updating eps in `mupdog`.

### Author(s)

David Gerard



---

obj\_for\_mu\_sigma2      *Objective function when updating mu and sigma2.*

---

**Description**

Objective function when updating mu and sigma2.

**Usage**

```
obj_for_mu_sigma2(mu, sigma2, phifk_mat, cor_inv, log_bb_dense)
```

**Arguments**

mu	A vector, the <i>i</i> th element is the variational posterior mean of individual <i>i</i> for the SNP.
sigma2	A vector, the <i>i</i> th element is the variational posterior variance of individual <i>i</i> for the SNP.
phifk_mat	A matrix that contains the standard normal quantile of the beta-binomial cdf at dosage <i>k</i> for individual <i>i</i> . The rows index the individuals and the columns index the dosages.
cor_inv	The inverse of the underlying correlation matrix.
log_bb_dense	A matrix of log-densities of the beta binomial. The rows index the individuals and the columns index the allele dosage. Allele dosage goes from -1 to ploidy, so there are ploidy + 2 elements.

**Value**

A double. The objective when updating mu and sigma2 in [mupdog](#).

**Author(s)**

David Gerard

---

obj\_for\_mu\_sigma2\_wrapper  
*Wrapper for [obj\\_for\\_mu\\_sigma2](#) so that I can use it in optim.*

---

**Description**

Wrapper for [obj\\_for\\_mu\\_sigma2](#) so that I can use it in optim.

**Usage**

```
obj_for_mu_sigma2_wrapper(muSigma2, phifk_mat, cor_inv, log_bb_dense)
```

**Arguments**

muSigma2	A vector. The first half are mu and the second half are sigma2.
phifk_mat	A matrix that contains the standard normal quantile of the beta-binomial cdf at dosage k for individual i. The rows index the individuals and the columns index the dosages.
cor_inv	The inverse of the underlying correlation matrix.
log_bb_dense	A matrix of log-densities of the beta binomial. The rows index the individuals and the columns index the allele dosage. Allele dosage goes from -1 to ploidy, so there are ploidy + 2 elements.

**Value**

A double. The objective when updating mu and sigma2 in [mupdog](#).

**Author(s)**

David Gerard

---

obj_for_rho	<i>Objective function when updating a single inbreeding coefficient.</i>
-------------	--

---

**Description**

Objective function when updating a single inbreeding coefficient.

**Usage**

```
obj_for_rho(rho, mu, sigma2, alpha, log_bb_dense, ploidy)
```

**Arguments**

rho	The inbreeding coefficient for the individual.
mu	A vector of posterior means. The jth element is the posterior mean of SNP j for the individual.
sigma2	A vector of posterior variances. The jth element is the posterior variance of SNP j for the individual.
alpha	A vector of allele frequencies. The jth element is the allele frequency for SNP j.
log_bb_dense	A matrix of log posterior densities. The rows index the SNPs and the columns index the dosage.
ploidy	The ploidy of the species.

**Value**

A double. The objective when updating rho in [mupdog](#).

**Author(s)**

David Gerard

---

obj\_for\_weighted\_lbb    *Objective function for updating the beta-binomial genotype distribution when model = "bb" in [flex\\_update\\_pivec](#).*

---

**Description**

Objective function for updating the beta-binomial genotype distribution when model = "bb" in [flex\\_update\\_pivec](#).

**Usage**

```
obj_for_weighted_lbb(parvec, ploidy, weight_vec)
```

**Arguments**

parvec	A vector of length 2. The first term is the current mean of the underlying beta. The second term is the current overdispersion parameter.
ploidy	The ploidy of the species.
weight_vec	A vector of length ploidy + 1 that contains the weights for each component beta-binomial.

**Value**

A double. The objective when updating the beta-binomial genotype distribution in [mupdog](#).

**Author(s)**

David Gerard

---

obj\_for\_weighted\_lnorm    *Objective function for updating discrete normal genotype distribution when model = "normal" in [flex\\_update\\_pivec](#).*

---

**Description**

Objective function for updating discrete normal genotype distribution when model = "normal" in [flex\\_update\\_pivec](#).

**Usage**

```
obj_for_weighted_lnorm(parvec, ploidy, weight_vec)
```

**Arguments**

parvec	A vector of length 2. The first term is the current mean of the underlying normal. The second term is the current standard deviation (not variance) of the normal.
ploidy	The ploidy of the species.
weight_vec	A vector of length ploidy + 1 that contains the weights for each component beta-binomial.

**Value**

A double. The objective when updating the normal genotype distribution in [mupdog](#).

**Author(s)**

David Gerard

---

oracle_cor	<i>Calculates the correlation between the true genotype and an oracle estimator.</i>
------------	--

---

**Description**

Calculates the correlation between the oracle MAP estimator (where we have perfect knowledge about the data generation process) and the true genotype. This is a useful approximation when you have a lot of individuals.

**Usage**

```
oracle_cor(n, ploidy, seq, bias, od, dist)
```

**Arguments**

n	The read-depth.
ploidy	The ploidy of the individual.
seq	The sequencing error rate.
bias	The allele-bias.
od	The overdispersion parameter.
dist	The distribution of the alleles.

**Details**

To come up with `dist`, you need some additional assumptions. For example, if the population is in Hardy-Weinberg equilibrium and the allele frequency is  $\alpha$  then you could calculate `dist` using the R code: `dbinom(x = 0:ploidy, size = ploidy, prob = alpha)`. Alternatively, if you know the genotypes of the individual's two parents are, say, `ref_count1` and `ref_count2`, then you could use the [get\\_q\\_array](#) function from the `updog` package: `get_q_array(ploidy)[ref_count1 + 1, ref_count2 + 1, ]`.

**Value**

The Pearson correlation between the true genotype and the oracle estimator.

**Author(s)**

David Gerard

**References**

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

**Examples**

```
## Hardy-Weinberg population with allele-frequency of 0.75.
## Moderate bias and moderate overdispersion.
## See how correlation decreases as we
## increase the ploidy.
ploidy <- 2
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
oracle_cor(n = 100, ploidy = ploidy, seq = 0.001,
           bias = 0.7, od = 0.01, dist = dist)

ploidy <- 4
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
oracle_cor(n = 100, ploidy = ploidy, seq = 0.001,
           bias = 0.7, od = 0.01, dist = dist)

ploidy <- 6
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
oracle_cor(n = 100, ploidy = ploidy, seq = 0.001,
           bias = 0.7, od = 0.01, dist = dist)
```

---

oracle\_cor\_from\_joint *Calculate the correlation of the oracle estimator with the true genotype from the joint distribution matrix.*

---

**Description**

Calculates the correlation between the oracle MAP estimator (where we have perfect knowledge about the data generation process) and the true genotype. This is a useful approximation when you have a lot of individuals.

**Usage**

```
oracle_cor_from_joint(jd)
```

**Arguments**

jd                    A matrix of numerics. Element (i, j) is the probability of genotype i - 1 and estimated genotype j - 1. This is usually obtained from [oracle\\_joint](#).

**Value**

The Pearson correlation between the true genotype and the oracle estimator.

**Author(s)**

David Gerard

**References**

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

**See Also**

[oracle\\_joint](#) for getting jd. [oracle\\_cor](#) for not having to first calculate jd.

**Examples**

```
## Hardy-Weinberg population with allele-frequency of 0.75.
## Moderate bias and moderate overdispersion.
ploidy <- 6
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
jd <- oracle_joint(n = 100, ploidy = ploidy, seq = 0.001,
                  bias = 0.7, od = 0.01, dist = dist)
oracle_cor_from_joint(jd = jd)

## Compare to oracle_cor
oracle_cor(n = 100, ploidy = ploidy, seq = 0.001,
          bias = 0.7, od = 0.01, dist = dist)
```

---

oracle\_joint

*The joint probability of the genotype and the genotype estimate of an oracle estimator.*

---

**Description**

This returns the joint distribution of the true genotypes and an oracle estimator given perfect knowledge of the data generating process. This is a useful approximation when you have a lot of individuals.

**Usage**

```
oracle_joint(n, ploidy, seq, bias, od, dist)
```

**Arguments**

n	The read-depth.
ploidy	The ploidy of the individual.
seq	The sequencing error rate.
bias	The allele-bias.
od	The overdispersion parameter.
dist	The distribution of the alleles.

**Details**

To come up with `dist`, you need some additional assumptions. For example, if the population is in Hardy-Weinberg equilibrium and the allele frequency is  $\alpha$  then you could calculate `dist` using the R code: `dbinom(x = 0:ploidy, size = ploidy, prob = alpha)`. Alternatively, if you know the genotypes of the individual's two parents are, say, `ref_count1` and `ref_count2`, then you could use the `get_q_array` function from the `updog` package: `get_q_array(ploidy)[ref_count1 + 1, ref_count2 + 1, ]`.

See the Examples to see how to reconcile the output of `oracle_joint` with `oracle_mis` and `oracle_mis_vec`.

**Value**

A matrix. Element  $(i, j)$  is the joint probability of estimating the genotype to be  $i+1$  when the true genotype is  $j+1$ . That is, the estimated genotype indexes the rows and the true genotype indexes the columns. This is when using an oracle estimator.

**Author(s)**

David Gerard

**References**

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

**See Also**

[oracle\\_plot](#) For visualizing the joint distribution output from `oracle_joint`.

[oracle\\_mis\\_from\\_joint](#) For obtaining the same results as `oracle_mis` directly from the output of `oracle_joint`.

[oracle\\_mis\\_vec\\_from\\_joint](#) For obtaining the same results as `oracle_mis_vec` directly from the output of `oracle_joint`.

[oracle\\_cor\\_from\\_joint](#) For obtaining the same results as `oracle_cor` directly from the output of `oracle_joint`.

**Examples**

```
## Hardy-Weinberg population with allele-frequency of 0.75.
## Moderate bias and moderate overdispersion.
ploidy <- 4
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
jd <- oracle_joint(n = 100, ploidy = ploidy, seq = 0.001,
                  bias = 0.7, od = 0.01, dist = dist)
jd

## Get same output as oracle_mis this way:
1 - sum(diag(jd))
oracle_mis(n = 100, ploidy = ploidy, seq = 0.001,
           bias = 0.7, od = 0.01, dist = dist)

## Get same output as oracle_mis_vec this way:
1 - diag(sweep(x = jd, MARGIN = 2, STATS = colSums(jd), FUN = "/"))
oracle_mis_vec(n = 100, ploidy = ploidy, seq = 0.001,
               bias = 0.7, od = 0.01, dist = dist)
```

---

oracle\_mis

*Calculate oracle misclassification error rate.*


---

**Description**

Given perfect knowledge of the data generating parameters, `oracle_mis` calculates the misclassification error rate, where the error rate is taken over both the data generation process and the allele-distribution. This is an ideal level of the misclassification error rate and any real method will have a larger rate than this. This is a useful approximation when you have a lot of individuals.

**Usage**

```
oracle_mis(n, ploidy, seq, bias, od, dist)
```

**Arguments**

<code>n</code>	The read-depth.
<code>ploidy</code>	The ploidy of the individual.
<code>seq</code>	The sequencing error rate.
<code>bias</code>	The allele-bias.
<code>od</code>	The overdispersion parameter.
<code>dist</code>	The distribution of the alleles.



## Details

To come up with `dist`, you need some additional assumptions. For example, if the population is in Hardy-Weinberg equilibrium and the allele frequency is  $\alpha$  then you could calculate `dist` using the R code: `dbinom(x = 0:ploidy, size = ploidy, prob = alpha)`. Alternatively, if you know the genotypes of the individual's two parents are, say, `ref_count1` and `ref_count2`, then you could use the `get_q_array` function from the `updog` package: `get_q_array(ploidy)[ref_count1 + 1, ref_count2 + 1, ]`.

## Value

A double. The oracle misclassification error rate.

## Author(s)

David Gerard

## References

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

## Examples

```
## Hardy-Weinberg population with allele-frequency of 0.75.
## Moderate bias and moderate overdispersion.
## See how oracle misclassification error rates change as we
## increase the ploidy.
ploidy <- 2
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
oracle_mis(n = 100, ploidy = ploidy, seq = 0.001,
           bias = 0.7, od = 0.01, dist = dist)

ploidy <- 4
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
oracle_mis(n = 100, ploidy = ploidy, seq = 0.001,
           bias = 0.7, od = 0.01, dist = dist)

ploidy <- 6
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
oracle_mis(n = 100, ploidy = ploidy, seq = 0.001,
           bias = 0.7, od = 0.01, dist = dist)
```

---

`oracle_mis_from_joint` *Get the oracle misclassification error rate directly from the joint distribution of the genotype and the oracle estimator.*

---

**Description**

Get the oracle misclassification error rate directly from the joint distribution of the genotype and the oracle estimator.

**Usage**

```
oracle_mis_from_joint(jd)
```

**Arguments**

`jd` A matrix of numerics. Element (i, j) is the probability of genotype i - 1 and estimated genotype j - 1. This is usually obtained from [oracle\\_joint](#).

**Value**

A double. The oracle misclassification error rate.

**Author(s)**

David Gerard

**References**

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

**See Also**

[oracle\\_joint](#) for getting `jd`. [oracle\\_mis](#) for not having to first calculate `jd`.

**Examples**

```
## Hardy-Weinberg population with allele-frequency of 0.75.
## Moderate bias and moderate overdispersion.
ploidy <- 6
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
jd <- oracle_joint(n = 100, ploidy = ploidy, seq = 0.001,
                  bias = 0.7, od = 0.01, dist = dist)
oracle_mis_from_joint(jd = jd)

## Compare to oracle_cor
oracle_mis(n = 100, ploidy = ploidy, seq = 0.001,
           bias = 0.7, od = 0.01, dist = dist)
```

---

oracle_mis_vec	Returns the oracle misclassification rates for each genotype.
----------------	---

---

### Description

Given perfect knowledge of the data generating parameters, `oracle_mis_vec` calculates the misclassification error rate at each genotype. This differs from `oracle_mis` in that this will *not* average over the genotype distribution to get an overall misclassification error rate. That is, `oracle_mis_vec` returns a vector of misclassification error rates *conditional* on each genotype.

### Usage

```
oracle_mis_vec(n, ploidy, seq, bias, od, dist)
```

### Arguments

n	The read-depth.
ploidy	The ploidy of the individual.
seq	The sequencing error rate.
bias	The allele-bias.
od	The overdispersion parameter.
dist	The distribution of the alleles.

### Details

This is an ideal level of the misclassification error rate and any real method will have a larger rate than this. This is a useful approximation when you have a lot of individuals.

To come up with `dist`, you need some additional assumptions. For example, if the population is in Hardy-Weinberg equilibrium and the allele frequency is  $\alpha$  then you could calculate `dist` using the R code: `dbinom(x = 0:ploidy, size = ploidy, prob = alpha)`. Alternatively, if you know the genotypes of the individual's two parents are, say, `ref_count1` and `ref_count2`, then you could use the `get_q_array` function from the `updog` package: `get_q_array(ploidy)[ref_count1 + 1, ref_count2 + 1, ]`.

### Value

A vector of numerics. Element  $i$  is the oracle misclassification error rate when genotyping an individual with actual genotype  $i + 1$ .

### Author(s)

David Gerard

### References

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

**Examples**

```
## Hardy-Weinberg population with allele-frequency of 0.75.
## Moderate bias and moderate overdispersion.
ploidy <- 4
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
om <- oracle_mis_vec(n = 100, ploidy = ploidy, seq = 0.001,
                    bias = 0.7, od = 0.01, dist = dist)
om

## Get same output as oracle_mis this way:
sum(dist * om)
oracle_mis(n = 100, ploidy = ploidy, seq = 0.001,
          bias = 0.7, od = 0.01, dist = dist)
```

---

oracle\_mis\_vec\_from\_joint

*Get the oracle misclassification error rates (conditional on true genotype) directly from the joint distribution of the genotype and the oracle estimator.*

---

**Description**

Get the oracle misclassification error rates (conditional on true genotype) directly from the joint distribution of the genotype and the oracle estimator.

**Usage**

```
oracle_mis_vec_from_joint(jd)
```

**Arguments**

jd                    A matrix of numerics. Element (i, j) is the probability of genotype i - 1 and estimated genotype j - 1. This is usually obtained from [oracle\\_joint](#).

**Value**

A vector of numerics. Element i is the oracle misclassification error rate when genotyping an individual with actual genotype i + 1.

**Author(s)**

David Gerard

**References**

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

**See Also**

[oracle\\_joint](#) for getting jd. [oracle\\_mis\\_vec](#) for not having to first calculate jd.

**Examples**

```
## Hardy-Weinberg population with allele-frequency of 0.75.
## Moderate bias and moderate overdispersion.
ploidy <- 6
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
jd <- oracle_joint(n = 100, ploidy = ploidy, seq = 0.001,
                  bias = 0.7, od = 0.01, dist = dist)
oracle_mis_vec_from_joint(jd = jd)

## Compare to oracle_cor
oracle_mis_vec(n = 100, ploidy = ploidy, seq = 0.001,
              bias = 0.7, od = 0.01, dist = dist)
```

---

oracle\_plot

*Construct an oracle plot from the output of [oracle\\_joint](#).*

---

**Description**

After obtaining the joint distribution of the true genotype with the estimated genotype from the oracle estimator using [oracle\\_joint](#), you can use `oracle_plot` to visualize this joint distribution.

**Usage**

```
oracle_plot(jd)
```

**Arguments**

jd                    A matrix containing the joint distribution of the true genotype and the oracle estimator. Usually, this is obtained by a call from [oracle\\_joint](#).

**Value**

A `ggplot` object containing the oracle plot. The x-axis indexes the possible values of the estimated genotype. The y-axis indexes the possible values of the true genotype. The number in cell (i, j) is the probability that an individual will have true genotype i but is estimated to have genotype j. This is when using an oracle estimator. The cells are also color-coded by the size of the probability in each cell. At the top are listed the oracle misclassification error rate and the correlation of the true genotype with the estimated genotype. Both of these quantities may be derived from the joint distribution.

**Author(s)**

David Gerard

**References**

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

**See Also**

[oracle\\_joint](#) for obtaining jd.

**Examples**

```
ploidy <- 6
dist <- stats::dbinom(0:ploidy, ploidy, 0.75)
jd <- oracle_joint(n = 100, ploidy = ploidy, seq = 0.001,
                 bias = 0.7, od = 0.01, dist = dist)
pl <- oracle_plot(jd = jd)
print(pl)
```

---

pbetabinom_double	<i>The distribution function of the betabinomial. This is generally only advisable if q is relatively small.</i>
-------------------	--

---

**Description**

The distribution function of the betabinomial. This is generally only advisable if q is relatively small.

**Usage**

```
pbetabinom_double(q, size, mu, rho, log_p)
```

**Arguments**

q	A quantile.
size	The total number of draws.
mu	The mean of the beta.
rho	The overdispersion parameter of the beta.
log_p	A logical. Should return the log-probability TRUE or not FALSE?

**Value**

The tail-probability of the beta-binomial.

**Author(s)**

David Gerard

---

pen_bias	<i>Penalty on bias parameter.</i>
----------	-----------------------------------

---

**Description**

Penalty on bias parameter.

**Usage**

```
pen_bias(h, mu_h, sigma2_h)
```

**Arguments**

h	The current value of bias parameter. Must be greater than 0. A value of 1 means no bias.
mu_h	The prior mean of the log-bias parameter.
sigma2_h	The prior variance (not standard deviation) of the log-bias parameter. Set to Inf to return 0.

**Value**

A double. The default penalty on the allelic bias parameter.

**Author(s)**

David Gerard

---

pen_seq_error	<i>Penalty on sequencing error rate.</i>
---------------	--

---

**Description**

Penalty on sequencing error rate.

**Usage**

```
pen_seq_error(eps, mu_eps, sigma2_eps)
```

**Arguments**

eps	The current value of sequencing error rate. Must be between 0 and 1.
mu_eps	The prior mean of the logit sequencing error rate.
sigma2_eps	The prior variance (not standard deviation) of the logit sequencing error rate. Set this to Inf to return 0.

**Value**

A double. The default penalty on the sequencing error rate.

**Author(s)**

David Gerard

---

pivec\_from\_segmat     *Function to get the segregation probabilities from the distributions of each component and the weights of each component.*

---

**Description**

Function to get the segregation probabilities from the distributions of each component and the weights of each component.

**Usage**

```
pivec_from_segmat(p1segmat, p2segmat, p1weight, p2weight)
```

**Arguments**

p1segmat	The matrix of segregation probabilities for each component for parent 1. The rows index the components and the columns index the number of alleles to segregate.
p2segmat	The matrix of segregation probabilities for each component for parent 2. The rows index the components and the columns index the number of alleles to segregate.
p1weight	A vector of weights for each component (row) of p1segmat.
p2weight	A vector of weights for each component (row) of p2segmat.

**Value**

A vector. The  $i$ th element is the probability of segregating  $i+1$  total A alleles.

**Author(s)**

David Gerard

**See Also**

This is mostly used in [update\\_pp\\_f1](#).



---

plot.flexdog                      Draw a genotype plot from the output of flexdog.

---

### Description

A wrapper for `plot_geno`. This will create a genotype plot for a single SNP.

### Usage

```
## S3 method for class 'flexdog'  
plot(x, use_colorblind = TRUE, ...)
```

### Arguments

x	A flexdog object.
use_colorblind	Should we use a colorblind-safe palette (TRUE) or not (FALSE)? TRUE is only allowed if the ploidy is less than or equal to 6.
...	Not used.

### Details

On a genotype plot, the x-axis contains the counts of the non-reference allele and the y-axis contains the counts of the reference allele. The dashed lines are the expected counts (both reference and alternative) given the sequencing error rate and the allele-bias. The plots are color-coded by the maximum-a-posterior genotypes. Transparency is proportional to the maximum posterior probability for an individual's genotype. Thus, we are less certain of the genotype of more transparent individuals. These types of plots are used in Gerard et. al. (2018) and Gerard and Ferrão (2019).

### Value

A `ggplot` object for the genotype plot.

### Author(s)

David Gerard

### References

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).
- Gerard, D. and Ferrão, L. F. V. (2019). Priors for Genotyping Polyploids. *bioRxiv*. doi: [10.1101/751784](https://doi.org/10.1101/751784).

### See Also

`plot_geno` The underlying plotting function.  
`flexdog` Creates a flexdog object.

plot.mupdog

*Draw a genotype plot from the output of [mupdog](#).*

---

**Description**

A wrapper for [plot\\_geno](#). This will create a genotype plot for a single SNP.

**Usage**

```
## S3 method for class 'mupdog'  
plot(x, index, use_colorblind = TRUE, ...)
```

**Arguments**

x	A mupdog object.
index	The column number of the gene to plot.
use_colorblind	Should we use a colorblind-safe palette (TRUE) or not (FALSE)? TRUE is only allowed if the ploidy is less than or equal to 6.
...	Not used.

**Details**

On a genotype plot, the x-axis contains the counts of the non-reference allele and the y-axis contains the counts of the reference allele. The dashed lines are the expected counts (both reference and alternative) given the sequencing error rate and the allele-bias. The plots are color-coded by the maximum-a-posterior genotypes. Transparency is proportional to the maximum posterior probability for an individual's genotype. Thus, we are less certain of the genotype of more transparent individuals. These types of plots are used in Gerard et. al. (2018) and Gerard and Ferrão (2019).

**Value**

A [ggplot](#) object for the genotype plot.

**Author(s)**

David Gerard

**References**

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).
- Gerard, D. and Ferrão, L. F. V. (2019). Priors for Genotyping Polyploids. *bioRxiv*. doi: [10.1101/751784](https://doi.org/10.1101/751784).

**See Also**

[plot\\_genotype](#) The underlying plotting function.  
[mupdog](#) Creates a mupdog object.

**Examples**

```
data("mupout")
plot(mupout, 4)
```

---

plot_genotype	<i>Make a genotype plot.</i>
---------------	------------------------------

---

**Description**

The x-axis is the counts of the non-reference allele, and the y-axis is the counts of the reference allele. Transparency is controlled by the `maxpostprob` vector. These types of plots are used in Gerard et. al. (2018) and Gerard and Ferrão (2019).

**Usage**

```
plot_genotype(refvec, sizevec, ploidy, p1ref = NULL, p1size = NULL,
              p2ref = NULL, p2size = NULL, geno = NULL, seq = 0, bias = 1,
              maxpostprob = NULL, p1geno = NULL, p2geno = NULL,
              use_colorblind = TRUE)
```

**Arguments**

<code>refvec</code>	A vector of non-negative integers. The number of reference reads observed in the individuals
<code>sizevec</code>	A vector of positive integers. The total number of reads in the individuals.
<code>ploidy</code>	A non-negative integer. The ploidy of the species.
<code>p1ref</code>	A vector of non-negative integers. The number of reference reads observed in parent 1 (if the individuals are all siblings).
<code>p1size</code>	A vector of positive integers. The total number of reads in parent 1 (if the individuals are all siblings).
<code>p2ref</code>	A vector of non-negative integers. The number of reference reads observed in parent 2 (if the individuals are all siblings).
<code>p2size</code>	A vector of positive integers. The total number of reads in parent 2 (if the individuals are all siblings).
<code>geno</code>	The individual genotypes.
<code>seq</code>	The sequencing error rate.
<code>bias</code>	The bias parameter.
<code>maxpostprob</code>	A vector of the posterior probabilities of being at the modal genotype.

p1geno Parent 1's genotype.  
 p2geno Parent 2's genotype.  
 use\_colorblind A logical. Should we use a colorblind safe palette (TRUE), or not (FALSE)? Only allowed if ploidy  $\leq$  6.

### Details

If parental genotypes are provided (p1geno and p2geno) then they will be colored the same as the offspring. Since they are often hard to see, a small black dot will also indicate their position.

### Value

A `ggplot` object for the genotype plot.

### Author(s)

David Gerard

### References

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).
- Gerard, D. and Ferrão, L. F. V. (2019). Priors for Genotyping Polyploids. *bioRxiv*. doi: [10.1101/751784](https://doi.org/10.1101/751784).

### Examples

```
data("snpmat")
refvec <- snpmat$counts[snpmat$snp == "SNP1"]
sizevec <- snpmat$size[snpmat$snp == "SNP1"]
ploidy <- 6
plot_geno(refvec = refvec, sizevec = sizevec, ploidy = ploidy)
```

---

post_prob	<i>Variational posterior probability of having dosage A alleles when the ploidy is ploidy, the allele frequency is alpha, the individual-specific overdispersion parameter is rho, the variational mean is mu, and the variational variance is sigma2.</i>
-----------	--

---

### Description

Variational posterior probability of having dosage A alleles when the ploidy is ploidy, the allele frequency is alpha, the individual-specific overdispersion parameter is rho, the variational mean is mu, and the variational variance is sigma2.

**Usage**

```
post_prob(dosage, ploidy, mu, sigma2, alpha, rho)
```

**Arguments**

dosage	The number of A alleles.
ploidy	The ploidy of the individual.
mu	The variational mean.
sigma2	The variational variance (not standard deviation).
alpha	The allele frequency.
rho	The individual's overdispersion parameter.

**Value**

The posterior probability of having dosage A alleles.

**Author(s)**

David

---

pp_brent_obj	<i>Objective function when doing Brent's method in <a href="#">update_pp_f1</a> when one parent only has two mixing components.</i>
--------------	---

---

**Description**

Objective function when doing Brent's method in [update\\_pp\\_f1](#) when one parent only has two mixing components.

**Usage**

```
pp_brent_obj(firstmixweight, probmat, pvec, weight_vec, alpha)
```

**Arguments**

firstmixweight	The mixing weight of the first component.
probmat	The rows index the components and the columns index the segregation amount. Should only have two rows.
pvec	The distribution of the other parent.
weight_vec	The weights for each element.
alpha	The mixing weight on the uniform component.

**Value**

The objective value, as calculated by taking a convolution using [convolve\\_up](#) of the mixing distribution and pvec, then putting that probability distribution through [f1\\_obj](#).

**Author(s)**

David Gerard

---

qbetabinom_double	<i>The quantile function of the beta-binomial distribution parameterized by mean and overdispersion parameter.</i>
-------------------	--

---

**Description**

The quantile function of the beta-binomial distribution parameterized by mean and overdispersion parameter.

**Usage**

```
qbetabinom_double(p, size, mu, rho)
```

**Arguments**

p	The lower tail probability.
size	The total number of draws.
mu	The mean of the beta.
rho	The overdispersion parameter of the beta.

**Value**

The quantile of the beta-binomial.

**Author(s)**

David Gerard

---

rbetabinom_int	<i>One draw from the beta-binomial distribution parameterized by mean and overdispersion parameter.</i>
----------------	---

---

**Description**

One draw from the beta-binomial distribution parameterized by mean and overdispersion parameter.

**Usage**

```
rbetabinom_int(size, mu, rho)
```

**Arguments**

size	The total number of draws.
mu	The mean of the beta.
rho	The overdispersion parameter of the beta.

**Value**

A random sample from the beta-binomial.

**Author(s)**

David Gerard

---

rflexdog	<i>Simulate GBS data from the <a href="#">flexdog</a> likelihood.</i>
----------	---

---

**Description**

This will take a vector of genotypes and a vector of total read-counts, then generate a vector of reference counts. To get the genotypes, you could use [rgeno](#). The likelihood used to generate read-counts is described in detail in Gerard et. al. (2018).

**Usage**

```
rflexdog(sizevec, geno, ploidy, seq = 0.005, bias = 1, od = 0.001)
```

**Arguments**

sizevec	A vector of total read-counts for the individuals.
geno	A vector of genotypes for the individuals. I.e. the number of reference alleles each individual has.
ploidy	The ploidy of the species.
seq	The sequencing error rate.
bias	The bias parameter. $\Pr(\text{a read after selected}) / \Pr(\text{A read after selected})$ .
od	The overdispersion parameter. See the Details of the rho variable in <a href="#">betabinom</a> .

**Value**

A vector the same length as sizevec. The ith element is the number of reference counts for individual i.

**Author(s)**

David Gerard

## References

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).
- Gerard, D. and Ferrão, L. F. V. (2019). Priors for Genotyping Polyploids. *bioRxiv*. doi: [10.1101/751784](https://doi.org/10.1101/751784).

## See Also

[rgeno](#) for a way to generate genotypes of individuals. [rbetabinom](#) for how we generate the read-counts.

## Examples

```
set.seed(1)
n      <- 100
ploidy <- 6

## Generate the genotypes of individuals from an F1 population,
## where the first parent has 1 copy of the reference allele
## and the second parent has two copies of the reference
## allele.
genovec <- rgeno(n = n, ploidy = ploidy, model = "f1",
                p1geno = 1, p2geno = 2)

## Get the total number of read-counts for each individual.
## Ideally, you would take this from real data as the total
## read-counts are definitely not Poisson.
sizevec <- stats::rpois(n = n, lambda = 200)

## Generate the counts of reads with the reference allele
## when there is a strong bias for the reference allele
## and there is no overdispersion.
refvec <- rflexdog(sizevec = sizevec, geno = genovec,
                  ploidy = ploidy, seq = 0.001,
                  bias = 0.5, od = 0)

## Plot the simulated data using plot_geno.
plot_geno(refvec = refvec, sizevec = sizevec,
          ploidy = ploidy, seq = 0.001, bias = 0.5)
```

---

rgeno

*Simulate individual genotypes from one of the supported flexdog models.*

---

## Description

This will simulate genotypes of a sample of individuals drawn from one of the populations supported by [flexdog](#). See the details of [flexdog](#) for the models allowed. These genotype distributions are described in detail in Gerard and Ferrão (2019).



**Usage**

```
rgeno(n, ploidy, model = c("hw", "bb", "norm", "ash", "f1", "s1", "f1pp",
  "s1pp", "flex", "uniform"), allele_freq = NULL, od = NULL,
  p1geno = NULL, p2geno = NULL, mode = NULL, pivec = NULL,
  mu = NULL, sigma = NULL, p1_pair_weights = NULL,
  p2_pair_weights = NULL)
```

**Arguments**

n	The number of observations.
ploidy	The ploidy of the species.
model	What form should the prior take? See Details in <a href="#">flexdog</a> .
allele_freq	If model = "hw", then this is the allele frequency of the population. For any other model, this should be NULL.
od	If model = "bb", then this is the overdispersion parameter of the beta-binomial distribution. See <a href="#">betabinom</a> for details. For any other model, this should be NULL.
p1geno	Either the first parent's genotype if model = "f1" (or model = "f1pp"), or the only parent's genotype if model = "s1" (or model = "s1pp"). For any other model, this should be NULL.
p2geno	The second parent's genotype if model = "f1" (or model = "f1pp"). For any other model, this should be NULL.
mode	If model = "ash", this is the center of the distribution. This should be a non-integer value (so the mode is either the floor or the ceiling of mode). For any other model, this should be NULL.
pivec	A vector of probabilities. If model = "ash", then this represents the mixing proportions of the discrete uniforms. If model = "flex", then element i is the probability of genotype i - 1. For any other model, this should be NULL.
mu	If model = "norm", this is the mean of the normal. For any other model, this should be NULL.
sigma	If model = "norm", this is the standard deviation of the normal. For any other model, this should be NULL.
p1_pair_weights	The mixing weights for the bivalent pairs output in <a href="#">get_bivalent_probs</a> at the lvec level of p1geno. If model = "f1pp" then this is for the first parent. If model = "s1pp", then this is for the only parent. This should be NULL for all other values of model.
p2_pair_weights	If model = "s1pp", these are the mixing weights for the bivalent pairs output in <a href="#">get_bivalent_probs</a> at the lvec level of p2geno for the second parent. This should be NULL for all other values of model.

**Details**

The allowable variable values of `allele_freq`, `od`, `p1geno`, `p2geno`, `pivec`, and `mode` varies based on the value of `model`. If `model = "ash"`, then only `mode` and `pivec` can be non-NULL. If `model = "flex"` then only `pivec` can be non-NULL. If `model = "hw"`, then only `allele_freq` can be non-NULL. If `model = "f1"` then only `p1geno` and `p2geno` can be non-NULL. If `model = "s1"` then only `p1geno` can be non-NULL. If `model = "uniform"`, then none of the above variables can be non-NULL. If `model = "bb"`, then only `allele_freq`, and `od` can be non-NULL. If `model == "norm"`, then only `mu` and `sigma` can be non-NULL.

**Value**

A vector of length `n` with the genotypes of the sampled individuals.

**Author(s)**

David Gerard

**References**

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).
- Gerard, D. and Ferrão, L. F. V. (2019). Priors for Genotyping Polyploids. *bioRxiv*. doi: [10.1101/751784](https://doi.org/10.1101/751784).

**Examples**

```
## F1 Population where parent 1 has 1 copy of the referenc allele
## and parent 2 has 4 copies of the reference allele.
ploidy <- 6
rgeno(n = 10, ploidy = ploidy, model = "f1", p1geno = 1, p2geno = 4)

## A population in Hardy-Weinberge equilibrium with an
## allele frequency of 0.75
rgeno(n = 10, ploidy = ploidy, model = "hw", allele_freq = 0.75)
```

---

snpdat

*GBS data from Shirasawa et al (2017)*

---

**Description**

Contains counts of reference alleles and total read counts from the GBS data of Shirasawa et al (2017) for the three SNP's used as examples in Gerard et. al. (2018).

**Usage**

snpdat

**Format**

A `tibble` with 419 rows and 4 columns:

**id** The identification label of the individuals.

**snp** The SNP label.

**counts** The number of read-counts that support the reference allele.

**size** The total number of read-counts at a given SNP.

**Value**

A `tibble`. See the Format Section.

**Source**

<https://doi.org/10.1038/srep44207>

**References**

- Shirasawa, Kenta, Masaru Tanaka, Yasuhiro Takahata, Daifu Ma, Qinghe Cao, Qingchang Liu, Hong Zhai, Sang-Soo Kwak, Jae Cheol Jeong, Ung-Han Yoon, Hyeong-Un Lee, Hideki Hirakawa, and Sahiko Isobe "A high-density SNP genetic map consisting of a complete set of homologous groups in autohexaploid sweetpotato (*Ipomoea batatas*)." *Scientific Reports* 7 (2017). DOI: 10.1038/srep44207
- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).

---

summary.mupdog

*Provides some summaries from the output of mupdog.*

---

**Description**

Returns mean-dosage for each individual at each SNP and the SNP-specific distribution of MAP dosages. The mean-dosage in particular might be of interest for downstream analyses.

**Usage**

```
## S3 method for class 'mupdog'
summary(object, ...)
```

**Arguments**

object            A mupdog object.  
 ...                Not used.

**Value**

A list with the following elements:

- **freq** A matrix with the frequency of the dosages (rows) for each SNP (columns).
- **mean\_dosage** The posterior mean dosage of an individual (rows) for each SNP (columns).

**Author(s)**

David Gerard

**Examples**

```
data(mupout)
msum <- summary(mupout)
msum$freq[, 1:5]
boxplot(msum$mean_dosage ~ mupout$map_dosage,
        xlab = "MAP Dosage", ylab = "Mean Dosage")
```

---

uitdewilligen

*Subset of individuals and SNPs from Uitdewilligen et al (2013).*

---

**Description**

A list containing a matrix of reference counts, a matrix of total counts, and the ploidy level (4) of the species. This is a subset of the data from Uitdewilligen et al (2013).

**Usage**

```
uitdewilligen
```

**Format**

A list containing three objects. Two matrices and a numeric scalar:

**refmat** A matrix of read counts containing the reference allele. The rows index the individuals and the columns index the SNPs.

**sizemat** A matrix of the total number of read counts. The rows index the individuals and the columns index the SNPs.

**ploidy** The ploidy level of the species (just 4).

**Value**

A list. See the Format Section.

**Source**

<https://doi.org/10.1371/journal.pone.0062355>

**References**

- Uitdewilligen, J. G., Wolters, A. M. A., Bjorn, B., Borm, T. J., Visser, R. G., & van Eck, H. J. (2013). [A next-generation sequencing method for genotyping-by-sequencing of highly heterozygous autotetraploid potato](#). *PLoS One*, 8(5), e62355.

**See Also**

[mupout](#): a mupdog fit of these data.

---

 uni\_em

---

*EM algorithm to fit weighted ash objective.*


---

**Description**

Solves the following optimization problem

$$\max_{\pi} \sum_k w_k \log\left(\sum_j \pi_j \ell_{jk}\right).$$

It does this using a weighted EM algorithm.

**Usage**

```
uni_em(weight_vec, lmat, pi_init, lambda, itermax, obj_tol)
```

**Arguments**

weight_vec	A vector of weights. Each element of weight_vec corresponds to a column of lmat.
lmat	A matrix of inner weights. The columns are the "individuals" and the rows are the "classes."
pi_init	The initial values of pivec. Each element of pi_init corresponds to a row of lmat.
lambda	The penalty on the pi's. Should be greater than 0 and really really small.
itermax	The maximum number of EM iterations to take.
obj_tol	The objective stopping criterion.

**Value**

A vector of numerics. The update of pivec in [flexdog\\_full](#).

**Author(s)**

David Gerard

---

uni_em_const	<i>EM algorithm to fit weighted ash objective with a uniform mixing component.</i>
--------------	--

---

**Description**

Solves the following optimization problem

$$\max_{\pi} \sum_k w_k \log(\alpha/(K+1) + (1-\alpha) \sum_j \pi_j \ell_{jk}).$$

It does this using a weighted EM algorithm.

**Usage**

```
uni_em_const(weight_vec, lmat, pi_init, alpha, lambda, itermx, obj_tol)
```

**Arguments**

weight_vec	A vector of weights. Each element of weight_vec corresponds to a column of lmat.
lmat	A matrix of inner weights. The columns are the "individuals" and the rows are the "classes."
pi_init	The initial values of pivec. Each element of pi_init corresponds to a row of lmat.
alpha	The mixing weight for the uniform component. This should be small (say, less than 10 <sup>-3</sup> ).
lambda	A vector of penalties on the pi's (corresponding to the rows of lmat). This can either be of length 1, in which case the same penalty is applied to each of the pi's. Or it can be the same length of pivec, in which case a different penalty is applied to each of the pi's. Larger penalties generally increase the value of the pi's, not shrink them.
itermx	The maximum number of EM iterations to take.
obj_tol	The objective stopping criterion.

**Value**

A vector of numerics. The update of pivec in [flexdog\\_full](#).

**Author(s)**

David Gerard

---

uni\_obj                      *Objective function optimized by [uni\\_em](#).*

---

**Description**

Objective function optimized by [uni\\_em](#).

**Usage**

```
uni_obj(pivec, weight_vec, lmat, lambda)
```

**Arguments**

pivec	The current parameters.
weight_vec	A vector of weights. Each element of weight_vec corresponds to a column of lmat.
lmat	A matrix of inner weights. The columns are the "individuals" and the rows are the "classes."
lambda	The penalty on the pi's. Should be greater than 0 and really really small.

**Value**

The objective optimized by [uni\\_em](#) during that separate unimodal EM algorithm.

**Author(s)**

David Gerard

---

uni\_obj\_const                      *Objective function optimized by [uni\\_em\\_const](#).*

---

**Description**

Objective function optimized by [uni\\_em\\_const](#).

**Usage**

```
uni_obj_const(pivec, alpha, weight_vec, lmat, lambda)
```

**Arguments**

pivec	The current parameters.
alpha	The mixing weight for the uniform component. This should be small (say, less than $10^{-3}$ ).
weight_vec	A vector of weights. Each element of weight_vec corresponds to a column of lmat.
lmat	A matrix of inner weights. The columns are the "individuals" and the rows are the "classes."
lambda	A vector of penalties on the pi's (corresponding to the rows of lmat). This can either be of length 1, in which case the same penalty is applied to each of the pi's. Or it can be the same length of pivec, in which case a different penalty is applied to each of the pi's. Larger penalties generally increase the value of the pi's, not shrink them.

**Value**

The objective optimized by [uni\\_em\\_const](#) during that separate unimodal EM algorithm.

**Author(s)**

David Gerard

---

update_dr	<i>Same as <a href="#">update_pp_f1</a> but I exclusively use the EM (instead of also Brent's method), and I allow for priors on the mixing proportions.</i>
-----------	--

---

**Description**

Same as [update\\_pp\\_f1](#) but I exclusively use the EM (instead of also Brent's method), and I allow for priors on the mixing proportions.

**Usage**

```
update_dr(weight_vec, model = c("f1pp", "f1ppdr"), control)
```

**Arguments**

weight_vec	colSums(wik_mat) from <a href="#">flexdog</a> . This is the sum of current posterior probabilities of each individual having genotype k.
model	The model to assume.
control	A list of anything else needed to be passed. E.g. if model = "ash", then inner_weights needs to be passed through control (see <a href="#">get_inner_weights</a> for how to get this matrix).



**Value**

A list with the following elements:

p1\_pair\_weights A list with the mixing weights for the bivalent components of parent 1.

p2\_pair\_weights A list with the mixing weights for the bivalent components of parent 2.

obj The maximized objective.

p1geno The estimated genotype for parent 1.

p2geno The estimated genotype for parent 2.

pivec The estimated genotype distribution for the offspring.

**See Also**

[update\\_pp\\_f1](#)

---

update_pp_f1	<i>Function to update the parameters in the preferential pairing F1 model.</i>
--------------	--

---

**Description**

Function to update the parameters in the preferential pairing F1 model.

**Usage**

```
update_pp_f1(weight_vec, control)
```

**Arguments**

weight\_vec colSums(wik\_mat) from [flexdog](#). This is the sum of current posterior probabilities of each individual having genotype k.

control A list of anything else needed to be passed. E.g. if model = "ash", then inner\_weights needs to be passed through control (see [get\\_inner\\_weights](#) for how to get this matrix).

**Value**

A list with the following elements:

p1\_pair\_weights A list with the mixing weights for the bivalent components of parent 1.

p2\_pair\_weights A list with the mixing weights for the bivalent components of parent 2.

obj The maximized objective.

p1geno The estimated genotype for parent 1.

p2geno The estimated genotype for parent 2.

pivec The estimated genotype distribution for the offspring.

**Author(s)**

David Gerard

**See Also**[update\\_dr](#)


---

update\_pp\_s1                      *Same as [update\\_pp\\_f1](#) but only allow s1.*

---

**Description**Same as [update\\_pp\\_f1](#) but only allow s1.**Usage**

update\_pp\_s1(weight\_vec, control)

**Arguments**

weight_vec	colSums(wik_mat) from <a href="#">flexdog</a> . This is the sum of current posterior probabilities of each individual having genotype k.
control	A list of anything else needed to be passed. E.g. if model = "ash", then inner_weights needs to be passed through control (see <a href="#">get_inner_weights</a> for how to get this matrix).

**Value**

A list with the following elements:

p1\_pair\_weights A list with the mixing weights for the bivalent components of parent 1.

p2\_pair\_weights A list with the mixing weights for the bivalent components of parent 2.

obj The maximized objective.

p1geno The estimated genotype for parent 1.

p2geno The estimated genotype for parent 2.

pivec The estimated genotype distribution for the offspring.

**See Also**[update\\_pp\\_f1](#)

---

update_R	<i>Update the underlying correlation matrix.</i>
----------	--

---

**Description**

Update the underlying correlation matrix.

**Usage**

```
update_R(postmean, postvar)
```

**Arguments**

postmean	The matrix of posterior means. The rows index the individuals and the columns index the SNPs.
postvar	The matrix of posterior variances. The rows index the individuals and the columns index the SNPs.

**Value**

A symmetric matrix of numerics. The update of the underlying correlation matrix.

**Author(s)**

David Gerard

---

updog	<i>updog Flexible Genotyping for Polyploids</i>
-------	---

---

**Description**

Implements empirical Bayes approaches to genotype polyploids from next generation sequencing data while accounting for allelic bias, overdispersion, and sequencing error. The main function is `flexdog`, which allows the specification of many different genotype distributions. An experimental function that takes into account varying levels of relatedness is implemented in `mupdog`. Also provided are functions to simulate genotypes (`rgeno`) and read-counts (`rflexdog`), as well as functions to calculate oracle genotyping error rates (`oracle_mis`) and correlation with the true genotypes (`oracle_cor`). These latter two functions are useful for read depth calculations. Run `browseVignettes(package = "updog")` in R for example usage. The methods are described in detail in Gerard et. al. (2018) and Gerard and Ferrão (2019).

## Details

The package is named updog for "Using Parental Data for Offspring Genotyping" because we originally developed the method for full-sib populations, but it works now for more general populations.

Our best competitor is probably the `fitPoly` package, which you can check out at <https://cran.r-project.org/package=fitPoly>. Though, we think that updog returns better calibrated measures of uncertainty when you have next-generation sequencing data.

If you find a bug or want an enhancement, please submit an issue at <http://github.com/dcgerard/updog/issues>.

## updog Functions

`flexdog` The main function that fits an empirical Bayes approach to genotype polyploids from next generation sequencing data.

`mupdog` An experimental approach to genotype autoployploids that accounts for varying levels of relatedness between the individuals in the sample.

`rgeno` simulate the genotypes of a sample from one of the models allowed in `flexdog`.

`rflexdog` Simulate read-counts from the `flexdog` model.

`plot.flexdog` Plotting the output of `flexdog`.

`plot.mupdog` Plotting the output of `mupdog`.

`oracle_joint` The joint distribution of the true genotype and an oracle estimator.

`oracle_plot` Visualize the output of `oracle_joint`.

`oracle_mis` The oracle misclassification error rate (Bayes rate).

`oracle_cor` Correlation between the true genotype and the oracle estimated genotype.

## updog Datasets

`snmdat` A small example dataset for using `flexdog`.

`uitdewilligen` A small example dataset for using `mupdog`.

`mupout` The output from fitting `mupdog` to `uitdewilligen`.

## Author(s)

David Gerard

## References

- Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. *Genetics*, 210(3), 789-807. doi: [10.1534/genetics.118.301468](https://doi.org/10.1534/genetics.118.301468).
- Gerard, D. and Ferrão, L. F. V. (2019). Priors for Genotyping Polyploids. *bioRxiv*. doi: [10.1101/751784](https://doi.org/10.1101/751784).

---

wem *Generalized version of [uni\\_em](#).*

---

### Description

Solves the following optimization problem

$$\max_{\pi} \sum_k w_k \log\left(\sum_j \pi_j \ell_{jk}\right).$$

It does this using a weighted EM algorithm.

### Usage

```
wem(weight_vec, lmat, pi_init, lambda, itermax, obj_tol)
```

### Arguments

weight_vec	A vector of weights. Each element of weight_vec corresponds to a column of lmat.
lmat	A matrix of inner weights. The columns are the "individuals" and the rows are the "classes."
pi_init	The initial values of pivec. Each element of pi_init corresponds to a row of lmat.
lambda	The penalty on the pi's. Should be greater than 0 and really really small.
itermax	The maximum number of EM iterations to take.
obj_tol	The objective stopping criterion.

### Value

A vector of numerics. The update of pivec in [flexdog\\_full](#).

### Author(s)

David Gerard

### See Also

[uni\\_em](#) for a description of the optimization problem.

**Examples**

```
set.seed(2)
n <- 3
p <- 5
lmat <- matrix(stats::runif(n * p), nrow = n)
weight_vec <- seq_len(p)
pi_init <- stats::runif(n)
pi_init <- pi_init / sum(pi_init)
wem(weight_vec = weight_vec,
     lmat       = lmat,
     pi_init    = pi_init,
     lambda     = 0,
     itermax    = 100,
     obj_tol    = 10^-6)
```

---

xi_double	<i>Adjusts allele dosage p by the sequencing error rate eps and the allele bias h.</i>
-----------	--

---

**Description**

Adjusts allele dosage p by the sequencing error rate eps and the allele bias h.

**Usage**

```
xi_double(p, eps, h)
```

**Arguments**

p	The allele dosage.
eps	The sequencing error rate.
h	The allele bias.

**Value**

The probability of a reference read adjusted by both the allele bias and the sequencing error rate.

**Author(s)**

David Gerard

---

xi_fun	<i>Adjusts allele dosage p by the sequencing error rate eps and the allele bias h.</i>
--------	--

---

**Description**

Adjusts allele dosage p by the sequencing error rate eps and the allele bias h.

**Usage**

```
xi_fun(p, eps, h)
```

**Arguments**

p	A vector of allele dosages.
eps	The sequencing error rate. Must either be of length 1 or the same length as p.
h	The allele bias. Must either be of length 1 or the same length as p.

**Value**

A vector of probabilities of the reference read adjusted by both the sequencing error rate and the allele bias.

**Author(s)**

David Gerard

# Index

## \*Topic **datasets**

- mupout, [54](#)
  - snpdatt, [82](#)
  - uitdewilligen, [84](#)
- ashpen\_fun, [4](#)
- betabinom, [25](#), [31](#), [79](#), [81](#)  
betabinom (dbetabinom), [8](#)
- compute\_all\_log\_bb, [5](#), [20](#)  
compute\_all\_phikf, [5](#)  
compute\_all\_post\_prob, [6](#), [20](#)  
convolve, [7](#)  
convolve\_up, [7](#), [77](#)
- dbernbinom, [8](#)  
dbetabinom, [8](#)  
dbetabinom\_alpha\_beta\_double, [10](#)  
dbetabinom\_double, [11](#), [13](#), [15](#)  
dc\_dtau, [11](#), [13](#), [15](#)  
df\_deps, [12](#)  
dlbeta\_dc, [12](#), [12](#), [15](#)  
dlbeta\_deps, [13](#)  
dlbeta\_dh, [14](#)  
dlbeta\_dtau, [12](#), [13](#), [14](#)  
dlbeta\_dxi, [15](#)  
doutdist, [16](#), [43](#)  
dpen\_deps, [16](#)  
dpen\_dh, [17](#)  
dr\_pen, [18](#)  
dxi\_df, [18](#)  
dxi\_dh, [19](#)
- elbo, [19](#)  
eta\_double, [20](#)  
eta\_fun, [21](#)  
expit, [21](#)
- f1\_obj, [22](#), [77](#)  
flex\_update\_pivec, [35](#), [59](#)
- flexdog, [4](#), [22](#), [27](#), [32–35](#), [40](#), [42](#), [43](#), [47](#), [48](#),  
[73](#), [79–81](#), [88–92](#)  
flexdog\_full, [22](#), [23](#), [25](#), [26](#), [27](#), [31](#), [33](#), [39](#),  
[42](#), [48](#), [85](#), [86](#), [93](#)  
flexdog\_obj, [33](#), [34](#)  
flexdog\_obj\_out, [34](#)
- get\_bivalent\_probs, [35](#), [36–38](#), [81](#)  
get\_bivalent\_probs\_dr, [36](#)  
get\_conv\_inner\_weights, [37](#)  
get\_dimname, [38](#)  
get\_hyper\_weights, [38](#)  
get\_inner\_weights, [35](#), [39](#), [88–90](#)  
get\_probk\_vec, [40](#), [41](#), [42](#)  
get\_q\_array, [40](#), [60](#), [63](#), [65](#), [67](#)  
get\_uni\_rep, [41](#)  
get\_wik\_mat, [42](#), [43](#)  
get\_wik\_mat\_out, [43](#)  
ggplot, [69](#), [73](#), [74](#), [76](#)  
grad\_for\_eps, [44](#), [44](#), [56](#)  
grad\_for\_mu\_sigma2, [45](#), [45](#)  
grad\_for\_mu\_sigma2\_wrapper, [45](#)  
grad\_for\_weighted\_lbb, [46](#)  
grad\_for\_weighted\_lnorm, [47](#)
- initialize\_pivec, [47](#)  
is.flexdog, [48](#)  
is.mupdog, [49](#)
- log\_sum\_exp, [50](#)  
log\_sum\_exp\_2, [50](#)  
logit, [49](#)
- mupdog, [20](#), [51](#), [54–60](#), [74](#), [75](#), [83](#), [91](#), [92](#)  
mupout, [54](#), [85](#), [92](#)
- obj\_for\_alpha, [55](#)  
obj\_for\_eps, [44](#), [56](#)  
obj\_for\_mu\_sigma2, [45](#), [46](#), [57](#), [57](#)  
obj\_for\_mu\_sigma2\_wrapper, [45](#), [57](#)  
obj\_for\_rho, [58](#)



obj\_for\_weighted\_lbb, [46](#), [59](#)  
obj\_for\_weighted\_lnorm, [47](#), [59](#)  
optim, [52](#)  
oracle\_cor, [26](#), [32](#), [60](#), [62](#), [63](#), [91](#), [92](#)  
oracle\_cor\_from\_joint, [61](#), [63](#)  
oracle\_joint, [62](#), [62](#), [66](#), [68–70](#), [92](#)  
oracle\_mis, [26](#), [32](#), [63](#), [64](#), [66](#), [67](#), [91](#), [92](#)  
oracle\_mis\_from\_joint, [63](#), [65](#)  
oracle\_mis\_vec, [63](#), [67](#), [69](#)  
oracle\_mis\_vec\_from\_joint, [63](#), [68](#)  
oracle\_plot, [63](#), [69](#), [92](#)

pbetabinom (dbetabinom), [8](#)  
pbetabinom\_double, [70](#)  
pen\_bias, [18](#), [71](#)  
pen\_seq\_error, [17](#), [71](#)  
pivec\_from\_segmat, [72](#)  
plot.flexdog, [26](#), [32](#), [73](#), [92](#)  
plot.mupdog, [55](#), [74](#), [92](#)  
plot\_geno, [24](#), [30](#), [73–75](#), [75](#)  
post\_prob, [6](#), [76](#)  
pp\_brent\_obj, [77](#)

qbetabinom (dbetabinom), [8](#)  
qbetabinom\_double, [78](#)

rbetabinom, [80](#)  
rbetabinom (dbetabinom), [8](#)  
rbetabinom\_int, [78](#)  
rflexdog, [26](#), [32](#), [79](#), [91](#), [92](#)  
rgeno, [26](#), [32](#), [79](#), [80](#), [80](#), [91](#), [92](#)

snmdat, [82](#), [92](#)  
summary.mupdog, [55](#), [83](#)

tibble, [83](#)

uitdewilligen, [54](#), [55](#), [84](#), [92](#)  
uni\_em, [85](#), [87](#), [93](#)  
uni\_em\_const, [37](#), [86](#), [87](#), [88](#)  
uni\_obj, [87](#)  
uni\_obj\_const, [87](#)  
update\_dr, [18](#), [88](#), [90](#)  
update\_pp\_f1, [37](#), [72](#), [77](#), [88](#), [89](#), [89](#), [90](#)  
update\_pp\_s1, [90](#)  
update\_R, [91](#)  
updog, [91](#)  
updog-package (updog), [91](#)

wem, [93](#)

xi\_double, [94](#)  
xi\_fun, [95](#)