

# *linemodels* R package

Matti Pirinen, University of Helsinki, Finland

Oct 5, 2022

**Summary** Estimation of effects of multiple explanatory variables on multiple outcome measures has become routine across life sciences with high-throughput molecular technologies. The `linemodels` R-package allows a probabilistic clustering of variables based on their observed effect sizes on two outcomes. User can specify each candidate model by three parameters:

- scale, i.e., the magnitude of the effects,
- slope, i.e., the multiplicative relationship between the expected values of the two effects, and
- correlation, i.e., the allowed deviation from the expectation.

`linemodels` provides functions to estimate model probabilities of each variable together with the proportion parameters of the models, by taking into account uncertainty in the effect estimates and a possible correlation of the two effect estimators. The package further allows for optimization of the model parameters and use of Gaussian mixture models as the prior distribution of effect sizes.

Available at: <https://github.com/mjpirinen/linemodels>

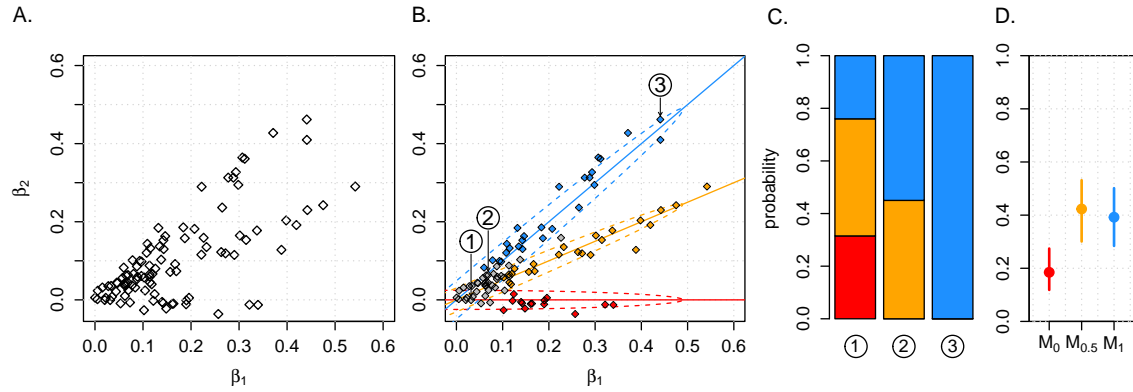
## 1. Motivating example

Consider that we have access to effect estimates of  $n = 100$  genetic variants on two subtypes  $Y_1$  and  $Y_2$  of a disease (Fig. 1A). Suppose that we are interested whether each variant is having

- (1) similar effects on both subtypes ( $\beta_1 \approx \beta_2$ ), or
- (2) a smaller effect for  $Y_2$  than for  $Y_1$  ( $\beta_2 < \beta_1$ ), or
- (3) no effect at all for  $Y_2$  ( $\beta_2 \approx 0$ ).

(While additional models might also be of interest, we keep this motivating example simple by only considering these three models.)

We apply `linemodels` by specifying the three models that are shown in Fig. 1B as colored lines. Output includes variable-specific posterior probabilities of the membership in each model (Fig. 1C) as well as proportion estimates for each model (Fig. 1D). From panel C we learn, for example, that while variant 3 has a very high probability of belonging to  $M_1$ , the data for variant 2 are indecisive between  $M_1$  and  $M_{0.5}$ . From panel D we see, for example, that about 18% (95% credible interval from 0.12 to 0.27) of the variants belong to  $M_0$ , and hence have  $\beta_2 \approx 0$ .



**Figure 1.** **A.** Effect estimates on  $n = 100$  variables on two outcomes. **B.** Three line models shown as colored lines with their 95% highest probability regions (dashed lines). The slopes of the models are 1 (model  $M_1$ , blue), 0.5 ( $M_{0.5}$ , orange) and 0 ( $M_0$ , red) corresponding to the above mentioned three relationships  $\beta_1 \approx \beta_2$ ,  $\beta_2 < \beta_1$  and  $\beta_2 \approx 0$ , respectively. Three variants are labelled as 1, 2 and 3. **C.** Variable-specific posterior probabilities of the membership in each model for the three variants labelled in panel **B.** **D.** Posterior distribution of proportion estimates of the models.

Note that the analysis using `linemodels` is not a trivial distance comparison between the points and the lines since

- we account for varying uncertainty of the effects on the two outcomes,
- we account for possible correlation between the effects on the two outcomes, e.g., due to shared control individuals in disease studies,
- we have defined Gaussian probability models around the lines.

With `linemodels` package, it is also possible to optimize the parameters of each line model rather than fix them before the analysis.

Next we show how to use the functions available in the `linemodels` package. Mathematical description of the model and algorithms can be found at the end of this document.

## 2. Install `linemodels`

To install `linemodels` from GitHub, you need to install ‘devtools’ package (in case you don’t have it).

```
install.packages("devtools")
```

After that, install `linemodels` with the following commands

```
library(devtools)
install_github("mjpirinen/linemodels")
library(linemodels)
```

If you can’t get the above working, all functions are in a single R file which you can also read in directly to R using command below

```
source("https://raw.githubusercontent.com/mjpirinen/linemodels/main/R/linemodels.R")
```

### 3. Specifying models

Each line model  $k$  is specified by **scale** ( $s_k$ ), **slope** ( $b_k$ ) and **correlation** ( $r_k$ ). We give these parameters as separate vectors **scales**, **slopes** and **cors**, each having  $K$  values where  $K$  is the number of models. Thus, to specify the following  $K = 3$  models from Fig.1,

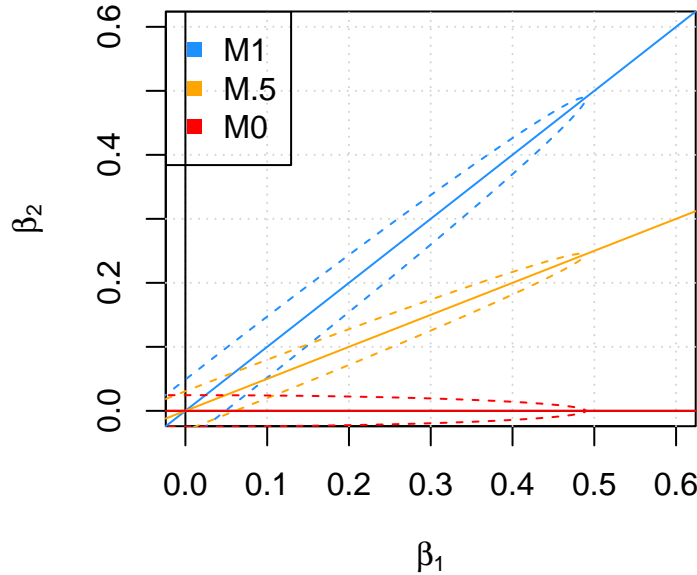
- $M_1$  with  $s = 0.2, b = 1, r = 0.995$ ,
- $M_{0.5}$  with  $s = 0.2, b = 0.5, r = 0.995$ ,
- $M_0$  with  $s = 0.2, b = 0, r = 0.995$ ,

we will define vectors:

```
scales = c(0.2, 0.2, 0.2)
slopes = c(1, 0.5, 0)
cors = c(0.995, 0.995, 0.995)
```

**visualize.line.models( )** We can visualize the models using **visualize.line.models()** to check that they seem reasonable. (Note: This function allows also other plotting parameters to be adjusted.)

```
visualize.line.models(scales, slopes, cors,
  model.names = c("M1", "M.5", "M0"),
  model.cols = c("dodgerblue", "orange", "red"),
  legend.position = "topleft",
  xlim = c(0, 0.6), ylim = c(0, 0.6),
  xlab = expression(beta[1]),
  ylab = expression(beta[2]))
```



**Figure 2.** The scale parameter  $s$  defines the expected range for the larger of the two effects as the scale is the standard deviation of the zero-centered effect size distribution. The correlation  $r$  defines how much deviation we allow around the exact line. Here the 95% regions determined by  $r$  are shown by the dashed lines.

## 4. Estimating line models

We will use the example data from Fig. 1 to demonstrate the analyses. The data is included in the package with name `linemodels_ex1` (and can also be found from GitHub page under `/data`).

```
dat = linemodels_ex1
head(dat, 2)
```

```
##          beta1          beta2          se1          se2
## 1 0.3394143 -0.012909064 0.01555013 0.01555013
## 2 0.1634416 -0.009596917 0.01510127 0.01510127
```

For each variable, we need its effect size estimates on the two outcomes (here `beta1` and `beta2`) and the corresponding standard errors of these estimates (`se1` and `se2`). Additionally, we will need the correlation `r.lkhood` between the estimators of the two effect sizes. Here we assume that the two estimates come from independent data sets, and hence `r.lkhood = 0`. If overlapping samples and/or correlated outcome measures were used, then a non-zero `r.lkhood` value should be specified. (See section 7 for examples.)

`line.models()` To estimate the model probabilities separately for each variant, we need to fix the prior probabilities of the models. Let's use equal priors.

```
model.priors = c(1/3, 1/3, 1/3)
model.names = c("M1", "M.5", "M0")
```

Now we can apply `line.models()` that gives us the probabilities of each variable in each model.

```
res = line.models(X = dat[,c("beta1", "beta2")], SE = dat[,c("se1", "se2")],
                  scales = scales, slopes = slopes, cors = cors,
                  model.names = model.names,
                  model.priors = model.priors,
                  r.lkhood = 0)
head(res, 2)
```

```
##          M1          M.5          M0
## [1,] 7.458412e-31 3.139269e-16 1.0000000
## [2,] 2.350610e-08 8.985146e-05 0.9999101
```

The result has one row per each variable and one column per each model. The example output above tells that, for both of the first two variables, the posterior probability is close to 1 that they follow the model  $M_0$  when the other two possibilities were  $M_1$  or  $M_{0.5}$ .

`line.models.with.proportions()` If we do not want to fix the prior probabilities of the line models before the analysis, we can assign a prior distribution on those probabilities and estimate them. The prior distribution is  $\text{Dirichlet}(a_1, \dots, a_K)$ , where user can give the values  $a_k$  in a  $K$ -dimensional vector called `diri.prior`. By default, we have  $a_k = \frac{1}{K}$  for all  $k \leq K$ . Since this approach uses an iterative Gibbs sampler, we can also specify the number of iterations (default 200) and the length of the burnin period (default 20). Larger values may improve accuracy in some data sets.

```
res.2 = line.models.with.proportions(
  X = dat[,c("beta1", "beta2")],
  SE = dat[,c("se1", "se2")],
  scales = scales, slopes = slopes, cors = cors,
  model.names = model.names,
  r.lkhood = 0,
  n.iter = 200, n.burnin = 20)
```

The result is a list with two members:

- **params** contain the posterior distribution of the proportion of variables belonging to each of the models (mean, 95% highest posterior region and standard deviation)
- **groups** has the posterior probabilities of each variable belonging to each of the models (similar to output of `line.models()` above)

```
res.2$params
```

```
##           mean    95%low    95%up      sd
## M1  0.3965014 0.2787163 0.5184582 0.06308448
## M.5 0.4128487 0.2977940 0.5308047 0.06135697
## M0   0.1906500 0.1162862 0.2641636 0.03733847
```

The proportion estimates above and their 95% intervals correspond to Fig.1D.

```
head(res.2$groups, 2)
```

```
##      M1 M.5 M0
## [1,]  0   0  1
## [2,]  0   0  1
```

In these data, the membership probabilities are highly similar between `line.models.with.proportions()` and `line.models()` since the estimated proportions of the three groups do not differ much from the equal prior assumption that we used in `line.models()`. Note that in `line.models.with.proportions()`, the unit of accuracy of the probabilities is determined by the number of iterations and, for example, with 200 iterations the unit is  $1/200 = 0.5\%$ . This explains why we see here probabilities that are exactly 1 whereas earlier with `line.models()` we saw some tiny deviations from 1, such as 0.9999101 for the second variable.

## 5. Optimizing parameters of line models

Sometimes we may not have a good reason to fix the parameters of a line model to any particular value of scale, slope or correlation, but rather we may want to estimate from the data a set of suitable values for the parameters.

**line.models.optimize( )** The optimization function takes in initial values of all parameters and a matrix of TRUE/FALSE entries named `par.include` that describes which parameters will be included in the optimization process. This matrix has one row and three columns per each model, where the columns correspond to scale, slope and correlation, respectively.

For example, let's optimize the slope of the second model that represents variants where  $\beta_1 > \beta_2 > 0$ . We may not want to assume that  $\beta_2 = 0.5 \cdot \beta_1$  in case that some other slope would clearly better describe the relationship in the data. Let's optimize over  $b_2$  and use an initial value of  $b_2 = 0.3$ .

```

#par.include here says that only the slope of 2nd model will be optimized
par.include = rbind(c(FALSE, FALSE, FALSE),
                    c(FALSE, TRUE, FALSE),
                    c(FALSE, FALSE, FALSE))
line.models.optimize(
  X = dat[,c("beta1", "beta2")],
  SE = dat[,c("se1", "se2")],
  par.include = par.include,
  init.scales = scales,
  init.slopes = c(0, 0.3, 1), #initial value for 2nd slope is 0.3
  init.cors = cors,
  model.priors = model.priors,
  model.names = model.names,
  r.lkhood = 0)

```

```

## Initial values
## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.3, 1
## cors: 0.995, 0.995, 0.995
## Optimizing w.r.t: slope2, over the range: (-Inf,Inf)
##
## iter: 1 ; log-lkhood: 254.16637
## proportions: 0.204, 0.339, 0.457
## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.5043, 1
## cors: 0.995, 0.995, 0.995
##
## iter: 2 ; log-lkhood: 254.46918
## proportions: 0.193, 0.384, 0.423
## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.5047, 1
## cors: 0.995, 0.995, 0.995
##
## iter: 3 ; log-lkhood: 254.57355
## proportions: 0.19, 0.404, 0.407
## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.5049, 1
## cors: 0.995, 0.995, 0.995
##
## iter: 4 ; log-lkhood: 254.61341
## proportions: 0.188, 0.412, 0.4
## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.5049, 1
## cors: 0.995, 0.995, 0.995
##
## iter: 5 ; log-lkhood: 254.62947
## proportions: 0.188, 0.415, 0.397
## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.505, 1
## cors: 0.995, 0.995, 0.995
##
## iter: 6 ; log-lkhood: 254.6361
## proportions: 0.187, 0.417, 0.396

```

```

## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.505, 1
## cors: 0.995, 0.995, 0.995
##
## iter: 7 ; log-likelihood: 254.63887
## proportions: 0.187, 0.417, 0.395
## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.505, 1
## cors: 0.995, 0.995, 0.995
##
## iter: 8 ; log-likelihood: 254.64002
## proportions: 0.187, 0.418, 0.395
## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.505, 1
## cors: 0.995, 0.995, 0.995
##
## iter: 9 ; log-likelihood: 254.64051
## proportions: 0.187, 0.418, 0.395
## scales: 0.2, 0.2, 0.2
## slopes: 0, 0.505, 1
## cors: 0.995, 0.995, 0.995

##      scales      slopes  cors
## [1,]    0.2 0.0000000 0.995
## [2,]    0.2 0.5049872 0.995
## [3,]    0.2 1.0000000 0.995

```

As the optimization proceeds, the function prints on the console the current value of all the parameters as well as the model log-likelihood that it aims to maximize. The final result shows that the EM-algorithm converged to value  $b_2 = 0.505$ . This is very close to the value 0.5 that we had been using, and hence updating this slope value in the line model analysis would not lead to a noticeable difference from the results of our earlier analyses.

**When to optimize?** We shouldn't just blindly optimize the model parameters for model comparison since we may overfit to the data and the optimized model may not anymore describe exactly that property which we wanted it to describe in the first place. In cases where we can come up with a particular model that exactly describes our hypothesis of interest, such as no effect on disease subtype 2 ( $\beta_2 = 0$ ) or same effect in both subtypes ( $\beta_1 = \beta_2$ ), then we should simply use those assumptions to define the models rather than using the optimization to improve the fit. On the other hand, optimization is useful when we can see from the data that there is a particular relationship that we want to model, but we do not know exactly what would be an appropriate parameter value to capture well that relationship.

## 6. A mixture of Gaussians as the prior distribution of effects

By default, the prior distribution of the larger effect is  $\mathcal{N}(0, s^2)$  where  $s$  is the scale parameter. (The larger effect is  $\beta_1$  when  $|b| \leq 1$  and  $\beta_2$  when  $|b| > 1$  where  $b$  is the slope parameter.) For the `line.models()` function, it is also possible to specify the effect size prior as a mixture of Gaussians. This can be done by defining `scale.weights` matrix, that has one row per each model and the number of columns is the number of mixture components ( $L$ ). Each row is normalized to sum to 1. Denote by  $w_{k\ell}$  the element  $(k, \ell)$  of `scale.weights` matrix. Then the prior for the larger effect in model  $k$  is

$$\sum_{\ell=1}^L w_{k\ell} \mathcal{N}\left(0, \left(\frac{s_k}{L - \ell + 1}\right)^2\right),$$

where  $s_k$  is the scale parameter of the model  $k$ . Note that  $s_k$  is given in standard deviation units and hence  $s_k^2$  is the corresponding variance. (In case the weights of each model are the same, we can also give the weights as a single vector instead of a matrix that had  $K$  identical rows.)

For example, if we want to model the effect size of models

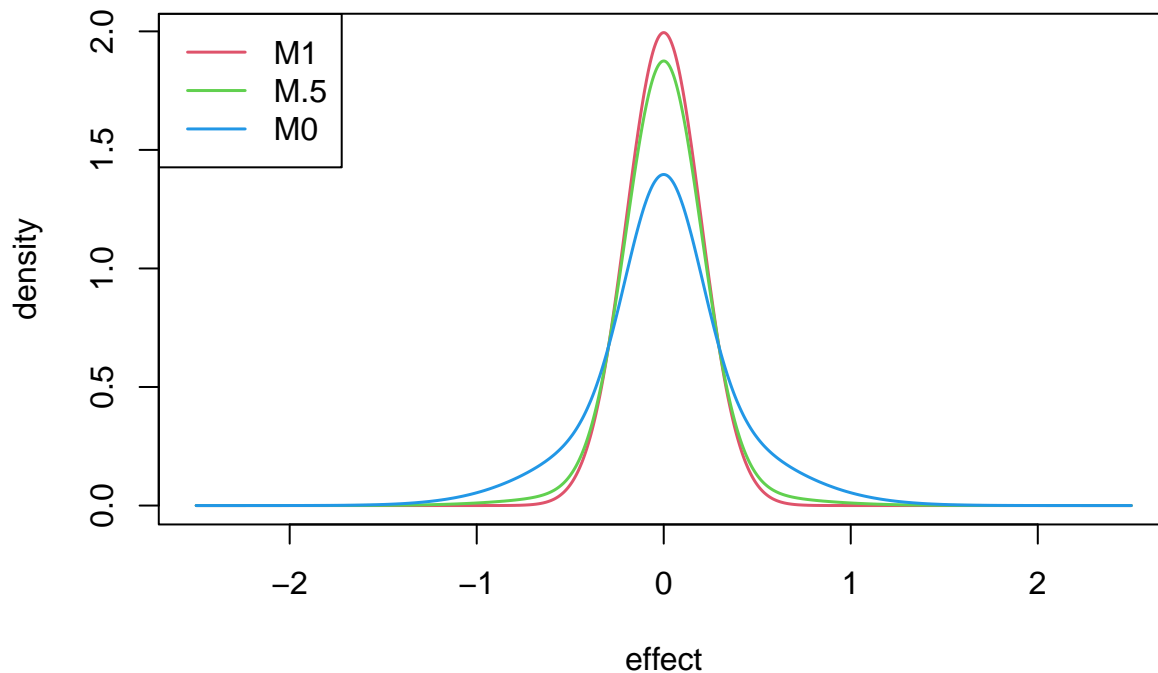
- $M_1$  as  $\mathcal{N}(0, 0.2^2)$ ,
- $M_{0.5}$  as a 10%:90% mixture of  $\mathcal{N}(0, 0.5^2)$  and  $\mathcal{N}(0, 0.2^2)$ ,
- $M_0$  as a 50%:50% mixture of  $\mathcal{N}(0, 0.5^2)$  and  $\mathcal{N}(0, 0.2^2)$ ,

then we can set all scales to  $s_k = 1.0$  and specify the `scale.weights` as

```
scales = c(1, 1, 1)
# elements of scale.weights will now corresp. to scales of
# 0.2 = 1.0/5, 0.25 = 1.0/4, 0.33 = 1.0/3, 0.5 = 1.0/2 and 1.0 = 1.0/1:
scale.weights = rbind(
  c(1.0, 0, 0, 0.0, 0.0), #N(0, 0.2^2)
  c(0.9, 0, 0, 0.1, 0.0), #0.9*N(0, 0.2^2) + 0.1*N(0, 0.5^2)
  c(0.5, 0, 0, 0.5, 0.0)) #0.5*N(0, 0.2^2) + 0.5*N(0, 0.5^2)
```

`visualize.scales( )` We can visualize the mixture prior distribution:

```
visualize.scales(scales, scale.weights, model.names)
```



```
## Plotted the following effect distributions:
## M1: Scale = 1 weights = 1,0,0,0,0.
## M.5: Scale = 1 weights = 0.9,0,0,0.1,0.
## M0: Scale = 1 weights = 0.5,0,0,0.5,0.
```

Applying these priors to effect sizes in `line.models()` is now straightforward: we just add `scale.weights` as an argument in the function call.



```
res.sc = line.models(X = dat[,c("beta1","beta2")], SE = dat[,c("se1","se2")],
                    scales = scales, slopes = slopes, cors = cors,
                    model.names = model.names,
                    model.priors = model.priors,
                    r.lkhood = 0,
                    scale.weights = scale.weights)
head(res.sc, 2)
```

```
##           M1           M.5           MO
## [1,] 7.573082e-31 2.160325e-07 0.9999998
## [2,] 3.455825e-08 1.840330e-03 0.9981596
```

## 7. How to set the value of `r.lkhood`?

When the effect estimators on the two outcomes are correlated, we must account that in our observation model through the parameter `r.lkhood`.

1. Suppose  $Y_1$  (outcome 1) and  $Y_2$  (outcome 2) are continuous outcomes and we have measured the effect of variable  $X$  on both of them separately on the same samples using linear regression models  $Y_1 = \mu_1 + X\beta_1 + \varepsilon_1$  and  $Y_2 = \mu_2 + X\beta_2 + \varepsilon_2$ . Then,

$$\text{r.lkhood} = \text{cor}(\hat{\beta}_1, \hat{\beta}_2) = \text{cor}(\varepsilon_1, \varepsilon_2) \approx \text{cor}(Y_1, Y_2),$$

where the last approximation assumes that  $X$  explains only a little of the variance of the outcome variables. Thus, to a first approximation, we may set `r.lkhood` equal to the correlation of the two outcome variables. If the samples are only partially overlapping, then we will need to shrink the correlation of the residuals towards zero by a factor of  $n_{12}/\sqrt{n_1 n_2}$ , where  $n_{12}$  is the number of samples that are shared between the analyses and  $n_1$  and  $n_2$  are the total sample sizes of each analysis.

2. If our data comes from two case-control analyses  $A$  and  $B$ , where samples are partially overlapping, we can derive a value for `r.lkhood` using the formula by Bhattacharjee et al. (2012). Denote by  $n_A^{(1)}$  the number of cases in  $A$  and by  $n_A^{(0)}$  the number of controls in  $A$ , and similarly for study  $B$ . Mark by  $n_{AB}^{(11)}$  and  $n_{AB}^{(00)}$ , respectively, the number of shared cases and shared controls between the studies and by  $n_{AB}^{(10)}$  and  $n_{AB}^{(01)}$ , respectively, the number of cases of  $A$  that are controls of  $B$  and the number of controls of  $A$  that are cases of  $B$ . Then

$$\text{r.lkhood} = \sqrt{\frac{n_A^{(1)} n_A^{(0)}}{n_A^{(1)} + n_A^{(0)}}} \sqrt{\frac{n_B^{(1)} n_B^{(0)}}{n_B^{(1)} + n_B^{(0)}}} \left( \frac{n_{AB}^{(11)}}{n_A^{(1)} n_B^{(1)}} - \frac{n_{AB}^{(10)}}{n_A^{(1)} n_B^{(0)}} - \frac{n_{AB}^{(01)}}{n_A^{(0)} n_B^{(1)}} + \frac{n_{AB}^{(00)}}{n_A^{(0)} n_B^{(0)}} \right).$$

## 8. Mathematical and technical details

Let  $\hat{\beta}_{ij}$  be the effect estimate of variable  $i = 1, \dots, n$  on outcome  $j = 1, 2$ , and  $\hat{\sigma}_{ij}$  its estimated standard error. Define, for  $k = 1, \dots, K$ , a line model  $\mathcal{M}_k$  via three parameters  $\mathcal{M}_k = (s_k, b_k, r_k)$  called scale  $s_k$ , slope  $b_k$  and correlation  $r_k$ , respectively. Intuitively,  $\mathcal{M}_k$  models the effects as centered around line  $\beta_{i2} = b_k \cdot \beta_{i1}$ , with larger of the two effects having prior standard deviation ("scale") of  $s_k$  and the deviation from the line determined by the correlation coefficient  $r_k$ . We first define the model for a diagonal case (slope = 1), and then use an orthogonal transformation to rotate the model to match the target slope.

For given  $\mathcal{M}_k = (s_k, b_k, r_k)$ , define the corresponding diagonal distribution of effect sizes as a bivariate Gaussian  $\mathcal{N}_2(\mathbf{0}, \mathbf{D}_k)$ , where the covariance matrix is

$$\mathbf{D}_k = \begin{pmatrix} 1 & r_k \\ r_k & 1 \end{pmatrix}.$$

Let  $\mathbf{T}_k$  be the rotation matrix that transforms the diagonal line to the line with slope  $b_k$ :

$$\mathbf{T}_k = \begin{pmatrix} \cos(\alpha) & -\sin(\alpha) \\ \sin(\alpha) & \cos(\alpha) \end{pmatrix}, \text{ where } \alpha = \arctan(b_k) - \frac{\pi}{4}.$$

The prior distribution of the effect sizes according to model  $\mathcal{M}_k$  is then defined as  $\mathcal{N}_2(\mathbf{0}, \boldsymbol{\Theta}_k)$ , where covariance matrix  $\boldsymbol{\Theta}_k = \frac{s_k^2}{m_k} \mathbf{T}_k \mathbf{D}_k \mathbf{T}_k^T$  and normalization by  $m_k = \max\{\mathbf{T}_k \mathbf{D}_k \mathbf{T}_k^T\}$  confirms that the larger of the standard deviations of the two effects is  $s_k$ . It is also possible to specify the prior distribution of the effect size as a mixture of Gaussians, e.g., to model heavier tails (see section 6 above).

In terms of the parameters defined above, we can use line models to define, for example,

- null model ( $s_k = 0$ ),
- independent effects model ( $s_k > 0, r_k = 0$ ),
- fixed effects model ( $s_k > 0, b_k = 1, r_k = 1$ ),
- outcome 1 specific effect ( $s_k > 0, b_k = 0, r_k = 1$ ),
- outcome 2 specific effect ( $s_k > 0, b_k = \infty, r_k = 1$ ),
- and any generalizations of these where deviation from exact line is allowed ( $0 \leq r_k < 1$ ).

The observation model for the effect size estimates is a Gaussian distribution around the true effect sizes with a covariance matrix

$$\hat{\boldsymbol{\Sigma}}_i = \begin{pmatrix} \hat{\sigma}_{i1}^2 & \hat{\sigma}_{i1}\hat{\sigma}_{i2}\hat{\rho}_i \\ \hat{\sigma}_{i1}\hat{\sigma}_{i2}\hat{\rho}_i & \hat{\sigma}_{i2}^2 \end{pmatrix},$$

where  $\hat{\rho}_i$  describes how the two effect size estimators are correlated, for example, because of the sample overlap in the data sets from where the two effects have been estimated.

It follows that by combining the Gaussian prior with the Gaussian observation model, the marginal distribution for the observed effect size estimates  $\hat{\boldsymbol{\beta}}_i = (\hat{\beta}_{i1}, \hat{\beta}_{i2})^T$  under model  $\mathcal{M}_k$  is

$$\hat{\boldsymbol{\beta}}_i | \mathcal{M}_k \sim \mathcal{N}_2(\mathbf{0}, \boldsymbol{\Theta}_k + \hat{\boldsymbol{\Sigma}}_i).$$

**8.1. Probabilistic assignment of variables into models** Given  $K$  line models  $(\mathcal{M}_k)_{k=1}^K$  and their prior probabilities  $\boldsymbol{\pi} = (\pi_k)_{k=1}^K$ , we can estimate the posterior probabilities that variable  $i$  follows the model  $k$  as

$$\Pr(i \sim \mathcal{M}_k | \text{Data}_i) = \frac{\pi_k \mathcal{N}_2(\hat{\boldsymbol{\beta}}_i | \mathbf{0}, \boldsymbol{\Theta}_k + \hat{\boldsymbol{\Sigma}}_i)}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}_2(\hat{\boldsymbol{\beta}}_i | \mathbf{0}, \boldsymbol{\Theta}_\ell + \hat{\boldsymbol{\Sigma}}_i)}.$$

This calculation can be done separately for each variable and is implemented in the `line.models()` function.

If we do not want to pre-specify the exact prior probabilities of each model, we can set a prior distribution on  $\boldsymbol{\pi}$  and estimate its posterior distribution together with assignment probabilities of variables into models. For this task, we use a prior distribution  $\boldsymbol{\pi} \sim \text{Dirichlet}(\delta_1, \dots, \delta_K)$ , where the default values of the hyperparameter are set to  $\delta_k = \frac{1}{K}$  for each  $k = 1, \dots, K$ . A Gibbs sampler to estimate the posterior distribution of this model is implemented in the `line.models.with.proportions()` function.

**8.2. EM-algorithm to optimize line model parameters** Using the `line.models.optimize()` function, one can optimize any set of model parameters (see section 5 above for an example).

Let's assume we have  $K$  line models and hence our model parameters are scales  $s_k$ , slopes  $b_k$ , correlations  $r_k$  and proportions  $\pi_k$  for each model  $k \leq K$ , but because of the constraint  $\sum_k \pi_k = 1$ , we can write  $\pi_K = 1 - \pi_1 - \dots - \pi_{K-1}$  and we have only  $4K - 1$  free parameters. Denote the  $3K$  line parameters by  $\boldsymbol{\Gamma} = \{s_k, b_k, r_k | k \leq K\}$  and the  $K$  proportion parameters by  $\boldsymbol{\pi} = \{\pi_k | k \leq K\}$ , and their values at iteration  $t$  of the algorithm as  $\boldsymbol{\Gamma}^{(t)}$  and  $\boldsymbol{\pi}^{(t)}$ , respectively, with the initial values corresponding to  $t = 0$ .

As is typical in EM-algorithms for mixture models, we introduce latent variables  $z_i \in \{1, \dots, K\}$  for each variable  $i = 1, \dots, n$  that tell to which line model the variable belongs. The augmented likelihood of the data consisting of the effect estimates  $(\hat{\beta}_i)_{i=1}^n$  and the corresponding standard errors  $(\hat{\sigma}_i)_{i=1}^n$  is

$$L(\mathbf{\Gamma}, \boldsymbol{\pi}, \mathbf{z}; \text{Data}) = p(\boldsymbol{\pi}) \prod_{i=1}^n p(z_i | \boldsymbol{\pi}) p(\text{Data}_i | z_i) = c \cdot \prod_{i=1}^n \pi_{z_i} \mathcal{N}_2 \left( \hat{\beta}_i | \mathbf{0}, \boldsymbol{\Theta}_{z_i} + \hat{\Sigma}_i \right),$$

where  $c$  is the normalizing constant of the  $K$ -dimensional uniform distribution of the proportion parameter  $\boldsymbol{\pi}$ . The log-likelihood is

$$\mathcal{L}(\mathbf{\Gamma}, \boldsymbol{\pi}, \mathbf{z}) = \log L(\mathbf{\Gamma}, \boldsymbol{\pi}, \mathbf{z}; \text{Data}) = \log(c) + \sum_{i=1}^n \left[ \log(\pi_{z_i}) + \log \mathcal{N}_2 \left( \hat{\beta}_i | \mathbf{0}, \boldsymbol{\Theta}_{z_i} + \hat{\Sigma}_i \right) \right].$$

The EM-algorithm at iteration  $t > 0$  consists of two steps.

1. Compute the posterior distribution  $(p_{ik}^{(t)})_{i,k}$  of latent variables  $(z_i)_{i=1}^n$  conditional on current values of parameters,  $\mathbf{\Gamma}^{(t-1)}$  and  $\boldsymbol{\pi}^{(t-1)}$ . This computation can be done separately for each variable  $i$  using the assignment probability formula

$$p_{ik}^{(t)} = \Pr(i \sim \mathcal{M}_k | \mathbf{\Gamma}^{(t-1)}, \boldsymbol{\pi}^{(t-1)}, \text{Data}_i) = \frac{\pi_k^{(t-1)} \mathcal{N}_2 \left( \hat{\beta}_i | \mathbf{0}, \boldsymbol{\Theta}_k^{(t-1)} + \hat{\Sigma}_i \right)}{\sum_{\ell=1}^K \pi_\ell^{(t-1)} \mathcal{N}_2 \left( \hat{\beta}_i | \mathbf{0}, \boldsymbol{\Theta}_\ell^{(t-1)} + \hat{\Sigma}_i \right)}.$$

2. Define new values  $\boldsymbol{\pi}^{(t)}$  and  $\mathbf{\Gamma}^{(t)}$  as maximizers of the expected log-likelihood function where the expectation is over the distribution of latent variables  $\mathbf{z}$ . Thus, the target function to be optimized is

$$\mathbb{E}_{\mathbf{z}^{(t)}} (\mathcal{L}(\mathbf{\Gamma}, \boldsymbol{\pi}, \mathbf{z})) = \log(c) + \sum_{i=1}^n \sum_{k=1}^K \left[ p_{ik}^{(t)} \log(\pi_k) + p_{ik}^{(t)} \log \mathcal{N}_2 \left( \hat{\beta}_i | \mathbf{0}, \boldsymbol{\Theta}_k + \hat{\Sigma}_i \right) \right].$$

We can analytically maximize with respect to  $\boldsymbol{\pi}$  by setting  $\pi_k^{(t)} = \frac{1}{n} \sum_i p_{ik}^{(t)}$ . For the remaining parameters in  $\mathbf{\Gamma}$ , we use numerical optimization of the above function.

The EM-algorithm is run until the difference in the maximized log-likelihood between adjacent iterations drops below a user given tolerance, such as the default value of  $10^{-3}$ .

## References

Bhattacharjee S., Rajaraman P., Jacobs K.B., et al. 2012. A Subset-Based Approach Improves Power and Interpretation for the Combined Analysis of Genetic Association Studies of Heterogeneous Traits, *The American Journal of Human Genetics*, 90(5): 821-835.