

# A simple guide of ENIGMA

Weixu Wang & Xiaolan Zhou

2022/1/20

## 1. Construct ENIGMA object

To build ENIGMA object, user is required to prepare bulk gene expression matrix (RNA-seq) and reference profile matrix, the reference profile matrix could be learned from scRNA-seq or FACS RNA-seq datasets. Here, I used example Non Small Cell Lung Cancer (NSCLC) single cell bulk RNA-seq datasets and single cell RNA-seq datasets to illustrate how to build ENIGMA object and perform deconvolution. The datasets could be downloaded from <https://doi.org/10.5281/zenodo.5906932>

```
# the example datasets could be downloaded from https://github.com/WXXkenmo/ENIGMA/tree/master
suppressPackageStartupMessages(library(Biobase))
suppressPackageStartupMessages(library(ENIGMA))
suppressPackageStartupMessages(library(Seurat))
suppressPackageStartupMessages(library(SingleCellExperiment))
dataNSCLC <- readRDS("/mnt/data1/weixu/HiDe/dataNSCLC.rds")
ref_sc <- readRDS("/mnt/data1/weixu/HiDe/ref.rds")

#We used the third patients to generate reference
ref_sc_sub <- ref_sc[,ref_sc$PatientID %in% "3" == TRUE]

## extract cells from tumor tissue.
ref_sc_sub <- ref_sc_sub[,ref_sc_sub$CellFromTumor %in% "1"]
ref_sc_sub <- ref_sc_sub[,ref_sc_sub$main_celltype %in%
                        c("Alveolar","Epi") == FALSE]

Bulk <- dataNSCLC[[5]]
Tumor <- dataNSCLC[[1]]
Immune <- dataNSCLC[[2]]
Endothelial <- dataNSCLC[[3]]
Fibroblast <- dataNSCLC[[4]]
pheno <- dataNSCLC[[6]]
# The pheno variable contain the label of each samples (LUSC vs LUAD)
names(pheno) <- colnames(Tumor)

###
ref_sc = exprs(ref_sc_sub)
metadata = pData(ref_sc_sub)
```

ENIGMA used three ways to construct the object

## 1. Using single cell RNA-seq datasets as reference

```
egm = create_ENIGMA(bulk = Bulk, ref = ref_sc, ref_type = "single_cell",  
                    meta_ref = metadata)
```

```
## Thu Jan 27 00:49:54 2022 Reference from Single Cell RNA-seq.
```

```
## Thu Jan 27 00:49:54 2022 Obtain reference from a matrix
```

The metadata contains the meta information of each cell in our reference datasets (which is required data.frame format). The reference profile (ref) need to be matrix.

```
head(metadata)
```

```
##           V1      V2 V3 dbCluster t_sub cell_type  
## AAACCTGAGTAACCCT_13 -27.09995 5.154716 NA B_cell <NA> B_cell  
## AAACCTGCAACGATGG_13 12.36870 4.525834 NA B_cell <NA> Myeloid  
## AAACCTGCATTAGCCA_13 20.98710 9.604408 NA B_cell <NA> Myeloid  
## AAACCTGGTCCATGAT_13 -19.22939 -7.118204 NA B_cell T_CD4 T_CD4  
## AAACCTGTCCTGTACC_13 27.25246 4.017953 NA B_cell <NA> Myeloid  
## AAACGGGCAAGAAAGG_13 30.32372 6.444399 NA B_cell <NA> Myeloid  
##           CellFromTumor      ClusterName PatientID main_celltype  
## AAACCTGAGTAACCCT_13      1 follicular B cells      3      B_cell  
## AAACCTGCAACGATGG_13      1      macrophages      3      Myeloid  
## AAACCTGCATTAGCCA_13      1      macrophages      3      Myeloid  
## AAACCTGGTCCATGAT_13      1      CD4+ T cells      3      T cell  
## AAACCTGTCCTGTACC_13      1      macrophages      3      Myeloid  
## AAACGGGCAAGAAAGG_13      1      macrophages      3      Myeloid
```

```
class(metadata)
```

```
## [1] "data.frame"
```

```
class(ref_sc)
```

```
## [1] "matrix" "array"
```

```
class(Bulk)
```

```
## [1] "matrix" "array"
```

ENIGMA also could adapt to the single cell object includes Seurat and SingleCellExperiment

```
## Create Seurat object for single cell profile
```

```
ref_se = CreateSeuratObject(ref_sc)
```

```
ref_se@meta.data = metadata
```

```
egm = create_ENIGMA(bulk = Bulk, ref = ref_se, ref_type = "single_cell")
```

```
## Thu Jan 27 00:49:55 2022 Reference from Single Cell RNA-seq.
```

```
## Thu Jan 27 00:49:55 2022 Obtain reference from a Seurat object
```

```
## The `data` slot in the default assay is used. The default assay is RNA
```

```
## The `meta.data` slot in the Seurat object is used as metadata of reference
```

```
## Create SingleCellExperiment for single cell profile
```

```
ref_se <- SingleCellExperiment(assays=list(counts = ref_sc))
```

```
colData(ref_se) = DataFrame(metadata)
```

```
egm = create_ENIGMA(bulk = Bulk, ref = ref_se, ref_type = "single_cell")
```

```
## Thu Jan 27 00:49:55 2022 Reference from Single Cell RNA-seq.
## Thu Jan 27 00:49:55 2022 Obtain reference from a SingleCellExperiment object
## The `counts` assay is used
## The `colData` assay in the SingleCellExperiment object is used as metadata of reference
```

## 2. Using FACS/Sorted RNA-seq to construct reference

User could also input FACS/Sorted RNA-seq as the reference, and ENIGMA would perform directly aggregation of FACS/Sorted RNA-seq according to cell type annotation to each samples

```
ref_m = cbind(Tumor,Immune,Endothelial,Fibroblast)
colnames(ref_m) = paste0("sample_",1:ncol(ref_m),sep="")
metadata = data.frame(celltype = c(rep("Tumor",ncol(Tumor)),rep("Immune",ncol(Immune)),
rep("Endothelial",ncol(Endothelial)),
rep("Fibroblast",ncol(Fibroblast))))
rownames(metadata) = colnames(ref_m)

# ref_type = "sort"
egm_sort = create_ENIGMA(bulk = Bulk, ref = ref_m, ref_type = "sort",meta_ref=metadata)
```

```
## Thu Jan 27 00:49:55 2022 Reference from FACS/Sort Bulk RNA-seq/microarray/scRNA-seq dataset.
## the metadata is required as data.frame format with merely one column of cell type annotation inform,
head(metadata)
```

```
##      celltype
## sample_1    Tumor
## sample_2    Tumor
## sample_3    Tumor
## sample_4    Tumor
## sample_5    Tumor
## sample_6    Tumor
```

## 3. Input reference profile matrix directly to build ENIGMA object

User could also input reference profile matrix directly to build ENIGMA object

```
ref_m = cbind(rowMeans(Tumor),rowMeans(Immune),
              rowMeans(Endothelial),rowMeans(Fibroblast))
colnames(ref_m) = c("Tumor","Immune","Endothelial","Fibroblast")

## ref_type = "aggre"
egm_aggre = ENIGMA::create_ENIGMA(bulk = Bulk, ref = ref_m, ref_type = "aggre")
```

```
## Thu Jan 27 00:49:55 2022 Reference from aggregated FACS/Sort Bulk RNA-seq/microarray/scRNA-seq.
```

### Note. count/normalized count inputs

We want to remind user that the input of bulk and ref need to be count or normalized count data, otherwise it would return warnings when construct ENIGMA object

```
egm_aggre = ENIGMA::create_ENIGMA(bulk = Bulk, ref = log2(ref_m+1), ref_type = "aggre")
```

```
## Thu Jan 27 00:49:55 2022 Reference from aggregated FACS/Sort Bulk RNA-seq/microarray/scRNA-seq.
```

```
## Warning in ENIGMA::create_ENIGMA(bulk = Bulk, ref = log2(ref_m + 1), ref_type
## = "aggre"): Reference matrix seems to be in log space. Please check if the
## reference matrix is normalized count data.
```

## 2. Cross-platform batch effects corrections

In order to remove the batch effects between reference and bulk, ENIGMA implement B/S-mode batch effects correction methods for users. When the `ref_type = "single_cell"`, ENIGMA would automatically use S-mode to correct batch effects. User is required to input the cluster id (e.g. macrostate or cell cluster, etc...) into `varname_cell_type`, and give a specific number of simulated pseudo-bulk samples (Default. 1000).

```
egm = batch_correct(egm, varname_cell_type = "main_celltype", n_pseudo_bulk=1000)
```

```
## Thu Jan 27 00:49:55 2022 Reference is from Single Cell RNA-seq, doing batch correction in S mode.
## Thu Jan 27 00:49:55 2022 Generating pseudo bulk...
## using 6 cores...
## Thu Jan 27 00:50:54 2022 Doing ComBat...
## Found 897 genes with uniform expression within a single batch (all zeros); these will not be adjusted
## Found2batches
## Adjusting for0covariate(s) or covariate level(s)
## Standardizing Data across genes
## Fitting L/S model and finding priors
## Finding parametric adjustments
## Adjusting the Data
## Thu Jan 27 00:51:04 2022 Restore reference...
```

When the `ref_type = "sort"` or `"aggre"`, ENIGMA would automatically use B-mode to correct batch effects. User doesn't need to input any parameter, just run

```
egm_sort = batch_correct(egm_sort)
```

```
## Thu Jan 27 00:51:11 2022 Reference is from FACS Bulk RNA-seq/microarray, doing batch correction in B
## Thu Jan 27 00:51:41 2022 Run B-mode to correct batch effect...
## Thu Jan 27 00:51:41 2022 Doing ComBat...
## Found 3633 genes with uniform expression within a single batch (all zeros); these will not be adjusted
## Found2batches
## Adjusting for0covariate(s) or covariate level(s)
## Standardizing Data across genes
## Fitting L/S model and finding priors
## Finding parametric adjustments
## Adjusting the Data
## Thu Jan 27 00:51:42 2022 Done.
```

If user doesn't want to perform batch effects correction, just skip this step.

### 3. Estimating cell type fraction matrix through reference based method

ENIGMA provide both robust linear regression (RLR) and CIBERSORT (CBS) methods to deconvolve bulk gene expression profile. The estimated fractions would be used in following CSE profile inference.

```
###Robust Linear Regression
```

```
egm = get_cell_proportion(egm, method = "RLR")
```

```
## Thu Jan 27 00:51:42 2022 Calculating cell type proportion of bulk samples...
```

```
## Using Robust Linear Regression...
```

```
###CIBERSORT
```

```
#egm = get_cell_proportion(egm, method = "CBS", nu.v = 0.5)
```

```
#Note: the CIBERSORT would cost for a long time to run.
```

### 4. Running ENIGMA

ENIGMA implement two type of model “trace norm” and “maximum L2 norm” (Default method). User could run maximum L2 norm model as follow

```
# random initialization
```

```
egm = ENIGMA_L2_max_norm(egm)
```

```
## Thu Jan 27 00:51:49 2022 Optimizing cell type specific expression profile...
```

```
## Converge in 9 steps
```

Also, user could run trace norm model as follow

```
# random initialization
```

```
egm = ENIGMA_trace_norm(egm,gamma = 1)
```

```
## Using admm solver...
```

```
## Normalization...
```

```
## Converge in 265 steps
```

ENIGMA also could track the trained model through model\_tracker.

```
# user also could specific the name of your trained model.
```

```
# If not, the model will be saved with the name as the format
```

```
# "model_type_the date you train the model".
```

```
# For example: maximum_L2_norm_model_Wed Jan 19 13:01:19 2022_trained
```

```
egm = ENIGMA_L2_max_norm(egm,alpha = 0.9,beta = 0.1,model_tracker=TRUE,
```

```
model_name = "model1",preprocess="log")
```

```
## Thu Jan 27 00:54:50 2022 Optimizing cell type specific expression profile...
```

```
## Converge in 20 steps
```

```
egm = ENIGMA_L2_max_norm(egm,alpha = 0.9,beta = 0.2,model_tracker=TRUE,
```

```
model_name = "model2",preprocess="log")
```

```
## Thu Jan 27 00:55:03 2022 Optimizing cell type specific expression profile...
```

```
## Converge in 12 steps
```

```
egm = ENIGMA_L2_max_norm(egm,alpha = 0.9,beta = 0.3,model_tracker=TRUE,
```

```
model_name = "model3",preprocess="log")
```

```
## Thu Jan 27 00:55:12 2022 Optimizing cell type specific expression profile...
```

```
## Converge in 9 steps
```

```
egm = ENIGMA_L2_max_norm(egm,alpha = 0.9,beta = 0.5,model_tracker=TRUE,
model_name = "Ken Wang",preprocess="log")
```

```
## Thu Jan 27 00:55:19 2022 Optimizing cell type specific expression profile...
## Converge in 6 steps
```

The list of trained models are saved in model\_name, according to the name of models, users could check the basic information of each model and their outputs

```
egm@model_name
```

```
##          Model Name          Model Type
## model1      model1 maximum L2 norm model
## model2      model2 maximum L2 norm model
## model3      model3 maximum L2 norm model
## model4      Ken Wang maximum L2 norm model
```

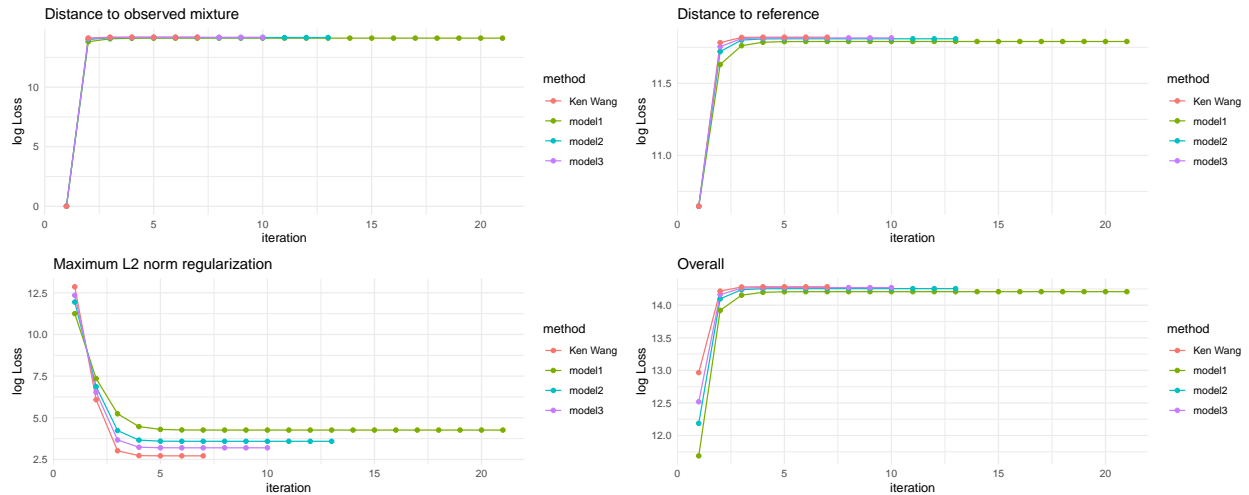
```
egm@model["Ken Wang"]
```

```
## $`Ken Wang`
## $`Ken Wang`$basic_infor
##   alpha beta step_size epsilon max_iter Normalize Normalize_method preprocess
## 1   0.9 0.5      0.01   0.001    1000      TRUE              frac        log
##   pos
## 1 TRUE
##
## $`Ken Wang`$result_CSE_normalized
## class: SingleCellExperiment
## dim: 17052 144
## metadata(0):
## assays(1): logcounts
## rownames(17052): A1BG A1CF ... ZZEF1 ZZZ3
## rowData names(0):
## colnames: NULL
## colData names(3): label sample cell_type
## reducedDimNames(0):
## altExpNames(0):
##
## $`Ken Wang`$loss_his
##           [,1]      [,2]      [,3]
## [1,] 1.020371e-26 42086.5 385391.98878
## [2,] 1.361415e+06 130900.0  436.74770
## [3,] 1.454620e+06 135714.8   19.43113
## [4,] 1.458980e+06 135925.9   14.30587
## [5,] 1.459172e+06 135934.3   14.10825
## [6,] 1.459181e+06 135934.5   14.09994
## [7,] 1.459181e+06 135934.5   14.09963
##
## $`Ken Wang`$result_CSE
## class: SingleCellExperiment
## dim: 17052 144
## metadata(0):
## assays(1): logcounts
## rownames(17052): A1BG A1CF ... ZZEF1 ZZZ3
## rowData names(0):
## colnames: NULL
```

```
## colData names(3): label sample cell_type
## reducedDimNames(0):
## altExpNames(0):
```

Users also could compare different models through inspecting their loss curves

```
plotLossCurve(egm,rlgType = "maximum L2 norm")
```



We could notice that the loss curve is not always decreasing, but this phenomena is relevant with the initialization point. Suppose we initialize the CSE as zero matrices

```
##clean all trained model
egm = clean_model(egm, NULL)
```

```
egm = ENIGMA::ENIGMA_L2_max_norm(egm,beta = 0.01,model_tracker=TRUE,
                                verbose=FALSE,model_name = "model1",
                                preprocess="log")
```

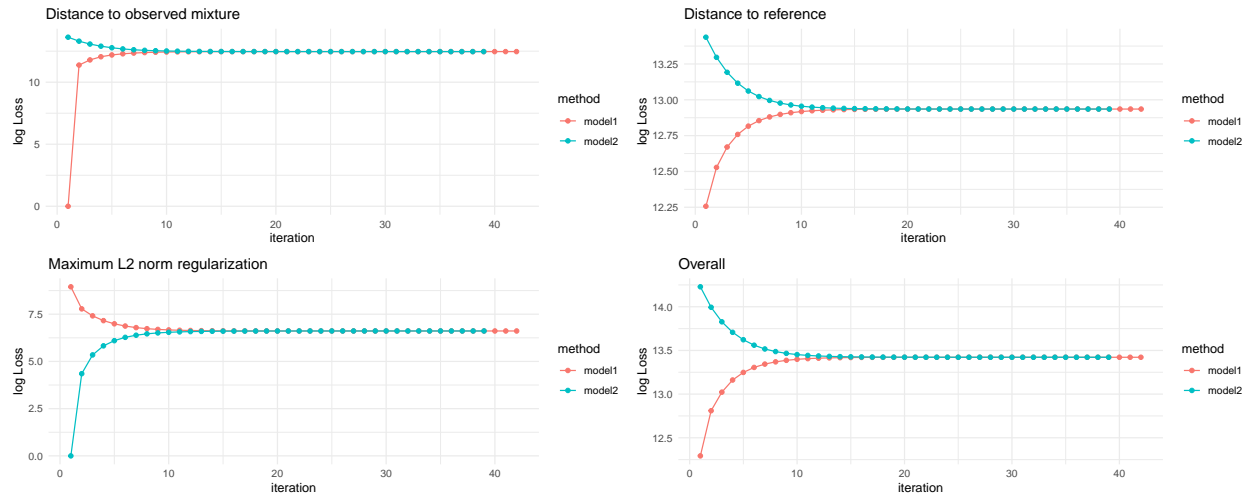
```
## Thu Jan 27 00:55:25 2022 Optimizing cell type specific expression profile...
## Converge in 41 steps
```

```
##
X_int = matrix(0, nrow = length(intersect(rownames(egm@bulk),rownames(egm@ref))),
              ncol = ncol(egm@bulk))
egm = ENIGMA::ENIGMA_L2_max_norm(egm,beta = 0.01,model_tracker=TRUE,
                                verbose=FALSE,model_name = "model2",
                                preprocess="log",
                                X_int = X_int)
```

```
## Thu Jan 27 00:55:50 2022 Optimizing cell type specific expression profile...
## Converge in 38 steps
```

Then, we could plot the loss curves as follow

```
plotLossCurve(egm,rlgType = "maximum L2 norm")
```



Therefore, users could use loss curve to see if their models' convergence. And all trained model would be saved in their ENIGMA object.