haoharryfeng / **ISLET**

<> **Code**   ⊙ Issues   ⨮ Pull requests   ▷ Actions   ⊞ Projects   ⊘ Security   ∿ Insights

**ISLET** / R / **dataprep.R** ⧉

haoharryfeng  0.99.6                                      8b8e137 · 2 years ago   ⟳ History

| Code | Blame | 87 lines (68 loc) · 3.43 KB | | Raw ⧉ ⬇ | ✎ ⌄ | <> |

```r
 1
 2    ###function to read-in and check data from case and control: observed expression, proportion, and sample-to-subject relationship
 3
 4  ∨ dataPrep<-function(dat_se){
 5      message("Begin: working on data preparation as the input for ISLET algorithm.")
 6    #  if (missing(case_dat_se) || missing(ctrl_dat_se))
 7    #    stop("SummarizedExperiment objects from both groups are needed.")
 8    #  if (length(intersect(colData(case_dat_se)[,1], colData(ctrl_dat_se)[,1]))>0)
 9    #    stop("Subject IDs across case group and control group must be unique.")
10    #  if (ncol(colData(case_dat_se))!=ncol(colData(ctrl_dat_se)))
11    #    stop("Case group and control group must have the same number of cell types.")
12
13      #subject id between cases and ctrls should also be unique, check and implement this later
14      #check for negative values, implement later
15
16      if (!is(dat_se, "SummarizedExperiment"))
17          stop("The input dataset must be a SummarizedExperiment object.")
18      if (length(unique(colData(dat_se)$group)) != 2)
19          stop("There must be two groups (case/ctrl) in the input SummarizedExperiment object.")
20      if (unique(colData(dat_se)$group)[1] != "case" || unique(colData(dat_se)$group)[2] != "ctrl")
21          stop("The names for the two groups in comparison should be labeled as
```

```
22              `case` and `ctrl` in the input SummarizedExperiment object.")
23
24
25     #separate cases and controls
26       idx <- which(colData(dat_se)$group == "case")
27
28       case_dat_se <- SummarizedExperiment(assays=list(counts=assays(dat_se)$counts[, idx]),
29                                           colData=colData(dat_se)[idx, -1])
30       ctrl_dat_se <- SummarizedExperiment(assays=list(counts=assays(dat_se)$counts[, -idx]),
31                                           colData=colData(dat_se)[-idx, -1])
32
33
34     #K = number of cell types
35       K <- ncol(colData(case_dat_se))-1
36
37     #N1 = number of samples for group 1
38       N1 <- ncol(assays(case_dat_se)$counts)
39     #N1 = number of samples for group 2
40       N2 <- ncol(assays(ctrl_dat_se)$counts)
41     #NS = total number of Samples for group 1&2
42       NS <- N1 + N2
43     #NU = total number of Unique subjects for group 1&2
44       caseUN <- length(unique(colData(case_dat_se)[, 1]))
45       ctrlUN <- length(unique(colData(ctrl_dat_se)[, 1]))
46       NU <- caseUN + ctrlUN
47
48
49
50       X_sub1 <- as.matrix(rbind(colData(case_dat_se)[, -1], colData(ctrl_dat_se)[, -1]))
51       X_sub2 <- rbind(matrix(1,  nrow=N1, ncol=K), matrix(0, nrow=N2, ncol=K))*X_sub1
52     #  X_sub2 = X_sub2[,1:para]
53       X_0 <- cbind(X_sub1, X_sub2)
54       X_list <- lapply(1, function(x){return(X_0)})
55       X <- bdiag(X_list)
56
```

Handwritten annotations:

Subject_ID    CT1  CT2  ...

Sample_ID 1                    1
Sample_ID 2                    1
...

$N_1$: # of case
$N_2$: # of ctrl

X_sub1: proportion $\in (0,1)$    $(N_1+N_2) \times K$

CT1 CT2...
sample1
sample2
:
sample N

X_sub2:
$\begin{bmatrix} // \end{bmatrix}$ 前 $N_1$行和X_sub1相同
                后 $N_2$行全是 0  $\in [0,1]$   $(N_1+N_2)\times K$

元素逐个对应相乘

X_0:

$$X = \begin{pmatrix}
\theta_{111} & \theta_{112} & \cdots & \theta_{11K} & z_1\theta_{111} & z_1\theta_{112} & \cdots & z_1\theta_{11K} \\
\theta_{121} & \theta_{122} & \cdots & \theta_{12K} & z_1\theta_{121} & z_1\theta_{122} & \cdots & z_1\theta_{12K} \\
\cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\
\theta_{1T_11} & \theta_{1T_12} & \cdots & \theta_{1T_1K} & z_1\theta_{1T_11} & z_1\theta_{1T_12} & \cdots & z_1\theta_{1T_1K} \\
\cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\
\theta_{J11} & \theta_{J12} & \cdots & \theta_{J1K} & z_J\theta_{J11} & z_J\theta_{J12} & \cdots & z_J\theta_{J1K} \\
\theta_{J21} & \theta_{J22} & \cdots & \theta_{J2K} & z_J\theta_{J21} & z_J\theta_{J22} & \cdots & z_J\theta_{J2K} \\
\cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\
\theta_{JT_J1} & \theta_{JT_J2} & \cdots & \theta_{JT_JK} & z_J\theta_{JT_J1} & z_J\theta_{JT_J2} & \cdots & z_J\theta_{JT_JK}
\end{pmatrix}_{N\times 2K}$$

$N = NS$ (# of sample)

K 列
$N_1$行 $\begin{cases} \end{cases}$ , $N_2$行 $\begin{cases} \end{cases}$
$\begin{bmatrix} 1 & 1 & 1 & 1 & \cdots \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$

Data inspector:
X    S4 [250 x 12] (Matrix::dgCMa  S4 object of class dgCMatrix
  i        integer [2250]          0 1 2 3 4 5 ...
  p        integer [13]            0 250 500 750 1000 1250 ...
  Dim      integer [2]             250 12
  Dimnames list [2]                List of length 2
  x        double [2250]           0.0223 0.0150 0.0131 0.0469 0.0469 0.0718 ...
  factors  list [0]                List of length 0

```
57      #obtain a vector of unique subject IDs, for all, to use later
58      sub_id <- c(colData(case_dat_se)[, 1], colData(ctrl_dat_se)[, 1])
59
60      propm <- as.matrix(rbind(colData(case_dat_se)[, -1], colData(ctrl_dat_se)[, -1]))
61    # propd = apply(propm, MARGIN = 2, makea, sub_id = sub_id, X = X, NU = NU, simplify = F)
62      propd <- apply(X=propm, MARGIN=2, FUN=makea,
63                     ind_id=sub_id, datX=X, aNU=NU, simplify=FALSE)
64
65      A_0 <- do.call(cbind, propd)
66      #A_list=lapply(1,function(x){return(A_0)})
67      A<-bdiag(A_0)
68
69      CT<-colnames(propm)
70
71      datuse <- inputSet(exp_case=assays(case_dat_se)$counts,
72                 exp_ctrl=assays(ctrl_dat_se)$counts,
73                 X=X,
74                 A=A,
75                 K=K,
76                 NS=NS,
77                 NU=NU,
78                 case_num=caseUN,
79                 ctrl_num=ctrlUN,
80                 CT=CT,
81                 SubjectID=sub_id,
82                 type='intercept'
83                 )
84      message("Complete: data preparation for ISLET.")
85      return(datuse)
86    }
```

Handwritten annotation near line 62: 3934



$$A = \begin{pmatrix} \mathbf{a}_{11} & 0 & 0 & 0 & \ldots & \mathbf{a}_{1K} & 0 & 0 & 0 \\ 0 & \mathbf{a}_{21} & 0 & 0 & \ldots & 0 & \mathbf{a}_{2K} & 0 & 0 \\ 0 & 0 & \ddots & 0 & \ldots & 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \mathbf{a}_{J1} & \ldots & 0 & 0 & 0 & \mathbf{a}_{JK} \end{pmatrix}_{N \times Q}$$

where $\mathbf{a}_{jk} := (\theta_{j_1 k}, \theta_{j_2 k}, \cdots, \theta_{j T_j k})'$ is simply a reorganized vector of cell type proportions, to align with random effect $\mathbf{u}$.　　(5)

Handwritten annotations (right margin):
$$A_{11} = \begin{bmatrix} \theta_{111} \\ \theta_{121} \\ \vdots \\ \theta_{1T_1 1} \end{bmatrix}$$

$N = NS = NU \times T$

第1个 subject 在第1个 cell type 的 different proportion across time

$Q = NU \times K$

$u \in \mathbb{R}^{Q \times 1}$

$Q = 50 \times 6 = 300$ in our case

$NS = 250$
$NU = 50$
$T = 5$

makea:

```
###function to make the design matrix [A] for random effect
#updated on 05/31/2022 to reflect the change in ID order
#user should sort their data by subject ID.
#makea <- function(onectprop, ind_id = sub_id, datX = X, aNU = NU){
makea <- function(onectprop, ind_id, datX , aNU){
  lp <- split(onectprop, ind_id)
  a1 <- matrix(0, nrow=nrow(datX), ncol=aNU)
#   ct_sub=table(sub_id)[as.character(unique(sub_id))]
#   lp=lp[names(ct_sub)]
  chk <- unique(ind_id) #chk should have the length of NU
  lp<-lp[as.character(chk)]
  count <- rep(0, length(chk))
  for(i in seq_len(aNU)){
    tmp <- sum(ind_id == chk[i])
    count[i] <- tmp
  }

  for(i in seq_len(aNU)){
    s <- 1+sum(count[0:(i-1)])
    e <- sum(count[seq_len(i)])
    a1[s:e, i] <- lp[[i]]
  }
}
  return(a1)
```

Handwritten annotations: "one cell type proportion" (above makea), "# of sample" (line 76), "# of subject" (line 77), "5个时间点" (near list), $0^{NS \times NU}$, "Gene1", "# of samples"



$$\mathbf{y} = X\boldsymbol{\beta} + A\mathbf{u} + \boldsymbol{\varepsilon}$$　　(3)

Handwritten: $N \times 1$, $y = \begin{Bmatrix} \vdots \\ N \end{Bmatrix}$ # of samples

where $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma_0^2 I)$ are the residuals. Here, $X$ and $A$ are the design matrices for the fixed-effect $\boldsymbol{\beta}$ and random-effect $\mathbf{u}$, respectively, where $\boldsymbol{\beta} = (m_1, m_2, \cdots, m_K, \beta_1, \beta_2, \cdots, \beta_K)'$ has two components: $(m_1, m_2, \cdots, m_K)$ are the baseline average gene expression in the control group, and $(\beta_1, \beta_2, \cdots, \beta_K)$ are the difference between the case group and the control group. The random effect $\mathbf{u} = (u_{11}, u_{21}, \cdots, u_{J1}, u_{12}, u_{22}, \cdots, u_{J2}, \cdots, u_{1K}, u_{2K}, \cdots, u_{JK})'$ captures the individual-level gene expression deviance from the group-level mean, for each cell type. The design matrices $X$ and $A$ are in the form:

Handwritten (bottom right): $R = 1 ( \frac{1}{2} - 4 cell type)$, $a_2$, $a_{12}$, $a_{11}$, $V1$, $V50$

| | V1 | V2 | V3 | V4 | |
|---|---|---|---|---|---|
| 1 | 0.02230079 | 0.00000000 | 0.00000000 | 0.00000000 | |
| 2 | 0.01501410 | 0.00000000 | 0.00000000 | 0.00000000 | |
| 3 | 0.01314183 | 0.00000000 | 0.00000000 | 0.00000000 | |
| 4 | 0.04691058 | 0.00000000 | 0.00000000 | 0.00000000 | |
| 5 | 0.04692466 | 0.00000000 | 0.00000000 | 0.00000000 | |
| 6 | 0.00000000 | 0.07183477 | 0.00000000 | 0.00000000 | |
| 7 | 0.00000000 | 0.01136901 | 0.00000000 | 0.00000000 | |
| 8 | 0.00000000 | 0.06145904 | 0.00000000 | 0.00000000 | |
| 9 | 0.00000000 | 0.02264138 | 0.00000000 | 0.00000000 | |
| 10 | 0.00000000 | 0.04052955 | 0.00000000 | 0.00000000 | |
| 11 | 0.00000000 | 0.00000000 | 0.03578315 | 0.00000000 | |
| 12 | 0.00000000 | 0.00000000 | 0.08965811 | 0.00000000 | |
| 13 | 0.00000000 | 0.00000000 | 0.05633649 | 0.00000000 | |
| 14 | 0.00000000 | 0.00000000 | 0.08530417 | 0.00000000 | |
| 15 | 0.00000000 | 0.00000000 | 0.01453683 | 0.00000000 | |
| 16 | 0.00000000 | 0.00000000 | 0.00000000 | 0.05460231 | |

☰  haoharryfeng / **ISLET**

🔍 Type `/` to search

<> **Code**    ⊙ Issues    ⅋ Pull requests    ▷ Actions    ⊞ Projects    ⊘ Security    📈 Insights

**ISLET** / R / **dataprep_slope.R** ⧉                                                      ⋯

🏞 **haoharryfeng** 0.99.6                                  8b8e137 · 2 years ago    🕐 History

| Code | Blame | 90 lines (73 loc) · 3.89 KB | | Raw ⧉ ⬇ | ✎ ⌄ | <> |

```r
1      ###function to read-in and check data from case and control: observed expression, proportion, and sample-to-subject relationship
2
3    ⌄ dataPrepSlope<-function(dat_se){
4          message("Begin: working on data preparation as the input for ISLET algorithm.")
5      #    if (missing(case_dat_se) || missing(ctrl_dat_se))
6      #        stop("SummarizedExperiment objects from both groups are needed.")
7      #    if (length(intersect(colData(case_dat_se)[,1], colData(ctrl_dat_se)[,1]))>0)
8      #        stop("Subject IDs across case group and control group must be unique.")
9      #    if (ncol(colData(case_dat_se))!=ncol(colData(ctrl_dat_se)))
10     #        stop("Case group and control group must have the same number of cell types.")
11
12         #subject id between cases and ctrls should also be unique, check and implement this later
13         #check for negative values, implement later
14         if (!is(dat_se, "SummarizedExperiment"))
15             stop("The input dataset must be a SummarizedExperiment object.")
16         if (length(unique(colData(dat_se)$group)) != 2)
17             stop("There must be two groups (case/ctrl) in the input SummarizedExperiment object.")
18         if (unique(colData(dat_se)$group)[1] != "case" || unique(colData(dat_se)$group)[2] != "ctrl")
19             stop("The names for the two groups in comparison should be
20                 labeled as `case` and `ctrl` in the input SummarizedExperiment object.")
21
```

```r
22          idx <- which(colData(dat_se)$group == "case")
23
24          case_dat_se <- SummarizedExperiment(assays=list(counts=assays(dat_se)$counts[, idx]),
25                                        colData=colData(dat_se)[idx, -1])
26          ctrl_dat_se <- SummarizedExperiment(assays=list(counts=assays(dat_se)$counts[, -idx]),
27                                        colData=colData(dat_se)[-idx, -1])
28
29
30
31
32          #K = number of cell types
33          K <- ncol(colData(case_dat_se))-2
34
35          #N1 = number of samples for group 1
36          N1 <- ncol(assays(case_dat_se)$counts)
37          #N1 = number of samples for group 2
38          N2 <- ncol(assays(ctrl_dat_se)$counts)
39          #NS = total number of Samples for group 1&2
40          NS <- N1 + N2
41          #NU = total number of Unique subjects for group 1&2
42          caseUN <- length(unique(colData(case_dat_se)[, 1]))
43          ctrlUN <- length(unique(colData(ctrl_dat_se)[, 1]))
44          NU <- caseUN + ctrlUN
45          case_age<-colData(case_dat_se)[, 2] ## The first two columns are: subject ID and sample age.
46          ctrl_age<-colData(ctrl_dat_se)[, 2]
47
48
49          X_sub1 <- as.matrix(rbind(colData(case_dat_se)[, -(seq_len(2))], colData(ctrl_dat_se)[, -(seq_len(2))]))
50          X_sub2 <- rbind(matrix(1, nrow=N1, ncol=K), matrix(0, nrow=N2, ncol=K))*X_sub1
51          X_age <- c(case_age, ctrl_age)
52          X_sub3 <- X_sub1*X_age
53          X_sub4 <- rbind(matrix(1, nrow=N1, ncol=K), matrix(0, nrow=N2, ncol=K))*X_sub1*X_age
54          ## This is difference in slope between two groups
55          X_sub4 <- X_sub4
56          X_0 <- cbind(X_sub1, X_sub2, X_sub3, X_sub4)
```

与 dataprep 基本相同. 此处 $\beta = (m_1, \cdots, m_K, \beta_1, \cdots, \beta_K,$

$\alpha_1, \cdots, \alpha_K, \gamma_1, \cdots, \gamma_K)^T \in \mathbb{R}^{4K}$

$$X_0 = \begin{pmatrix} \text{dataprep的 } X_0 & : & \theta_{\cdots} \cdot t_1 & : & \theta_{\cdots} \cdot t_1 \cdot Z_1 \\ & : & X\text{-sub3} & : & X\text{-sub4} \end{pmatrix} \in [0,1)^{N \times 4K}$$

$$\beta \in \mathbb{P}^{4k \times G}$$

```
57          X_list <- lapply(1, function(x){return(X_0)})
58          X <- bdiag(X_list)
59
60          #obtain a vector of unique subject IDs, for all, to use later
61          sub_id <- c(colData(case_dat_se)[, 1], colData(ctrl_dat_se)[, 1])
62
63          propm <- as.matrix(rbind(colData(case_dat_se)[, -(seq_len(2))], colData(ctrl_dat_se)[, -(seq_len(2))]))
64          # propd = apply(propm, MARGIN = 2, makea, sub_id = sub_id, X = X, NU = NU, simplify = F)
65          propd <- apply(X=propm, MARGIN=2, FUN=makea,
66                        ind_id=sub_id, datX=X, aNU=NU, simplify=FALSE)
67
68      A_0 <- do.call(cbind, propd)
69      #A_list=lapply(1,function(x){return(A_0)})
70      A<-bdiag(A_0)
71
72      CT<-colnames(propm)
73
74      datuse <- inputSet(exp_case=assays(case_dat_se)$counts,
75                        exp_ctrl=assays(ctrl_dat_se)$counts,
76                        X=X,
77                        A=A,
78                        K=K,
79                        NS=NS,
80                        NU=NU,
81                        case_num=caseUN,
82                        ctrl_num=ctrlUN,
83                        CT=CT,
84                        SubjectID=sub_id,
85                        type='slope'
86                        )
87      message("Complete: data preparation for ISLET.")
88      return(datuse)
89  }
```