

Algorithm SMMB

Convention: Sets of variables (*i.e.* ensembles) will be denoted in bold font.

Algorithm 1 Stochastic Multiple Markov blanket

SMMB(**r**, **t**, **X**, **T**, **K**, **k**, **m**, α): **consensus_MB**

INPUT:

r, maximal number of Markov blankets ($|\mathbf{MBs}|$) output by the algorithm
t, maximal number of iterations at top level,
X, observed data for genotypes, data matrix of dimension $n * p$, with:
 n, the number of individuals, and
 p, the number of variables
T, observed data for phenotype (vector of dimension *n*)
K, total number of variables sampled from **X** to learn a Markov blanket (top level)
k, number of variables sampled from **K** variables (inner level), $k < K$
m, maximal number of resamplings (of **k** variables) as long as the Markov blanket remains empty
 α , type I error threshold

OUTPUT:

consensus_MB, a consensus from all learned Markov blankets **MBs**, $|\mathbf{MBs}| \leq \mathbf{r}$

```

1: MBs  $\leftarrow \emptyset$ 
2: i  $\leftarrow 0$ 
3: while ( $|\mathbf{MBs}| \leq r$  and  $i \leq t$ )
4:   X*  $\leftarrow \text{sampling\_without\_replacement}(K, \mathbf{X})$ 
5:   MB*  $\leftarrow \text{learnMB}(\mathbf{X}^*, T, k, m, \alpha)$ 
6:   if not empty(MB*) then add(MBs, MB*) end if
7:   incr(i)
8: end while
9: consensus_MB  $\leftarrow \text{buildConsensus}(\mathbf{MBs}, \alpha)$ 
10: return consensus_MB

```

Algorithm 2 learnMB(**X***, **T**, **k**, **m**, α): **MB**

OUTPUT:

MB, a Markov blanket, possibly empty

```

1: MB  $\leftarrow \emptyset$  /*initialization of candidate Markov blanket*/
2: i  $\leftarrow 0$ 
3: repeat
4:   S  $\leftarrow \text{sampling\_without\_replacement}(k, \mathbf{X}^*)$ 

   /*Forward step*/
5:   s  $\leftarrow \text{argmax}_{\mathbf{s}' \subseteq \mathbf{S}} \{\text{assoc\_score}(\mathbf{s}', T, \mathbf{MB})\}$ 
6:   if not significant\_indepMB(s, T, MB,  $\alpha$ ) then
7:     MB  $\leftarrow \mathbf{MB} \cup \mathbf{s}$ 
   /*Backward step*/
8:   for each X  $\in \mathbf{MB}$ 
9:     for each S  $\subseteq \mathbf{MB} \setminus \{X\}$ , S  $\neq \emptyset$ 
10:      if(significant\_independence(X, T, S,  $\alpha$ )) then MB  $\leftarrow \mathbf{MB} \setminus \{X\}$ ; break end if
11:    end for
12:  end for
13: end if
14: incr(i)
15: until ((not empty(MB)) and (MB does not change)) or (empty(MB) and i = m)
16: return MB

```

Algorithm 3 buildConsensus(MBs, α): consensus_MB

INPUT:

MBs, a set of Markov blankets
 α , type I error threshold

OUTPUT:

consensus_MB, a consensus from all Markov blankets MBs

```
1: consensus_MB  $\leftarrow \bigcup_{\text{MB} \in \text{MBs}} \text{MB}$  /* initialization of consensus */  
   /*Backward step*/  
2: for each  $X \in \text{consensus\_MB}$   
3:   for each  $\text{S} \subseteq \text{consensus\_MB} \setminus \{X\}, \text{S} \neq \emptyset$   
4:     if (significant_independence( $X, T, \text{S}, \alpha$ )) then  
5:       consensus_MB  $\leftarrow \text{consensus\_MB} \setminus \{X\}$ ; break  
6:     end if  
7:   end for  
8: end for  
  
9: return consensus_MB
```

Algorithm 4 *assoc_score*($\text{S}_1, T, \text{S}_2$) : maximal score

INPUT:

S_1 , a set of variables
 T , a variable, $T \notin \text{S}_1$
 S_2 , a set of variables, $T \notin \text{S}_2$

```
1: score_max  $\leftarrow -\infty$   
2: for each  $X \in \text{S}_1$   
3:   stat  $\leftarrow \text{stat\_independence\_test}(X, T, \text{S}_2 \cup (\text{S}_1 \setminus \{X\}))$   
4:   if (stat > score_max) then score_max  $\leftarrow \text{stat}$ ; memorize(p-value) end if  
5: end for  
6: return score_max
```

Comments

Algorithms 2 (learnMB) and 3 (buildConsensus) use function *significant_independence*(X, T, S, α), with $X \neq T$ and $X \notin \text{S}$. Function *significant_independence*(X, T, S, α) runs a G-test of independence between variables X and T , conditional on set S . We denote stat_o the observed statistic returned by the test for *Stat*, the random variable with an unknown distribution. For the G-test, the distribution P_{H_0} of *Stat* under the hypothesis of independence H_0 is known, it is the *Chi – Squared* law. Thus, the function *significant_independence* returns true if and only if $P_{H_0}(\text{Stat} \geq \text{stat}_o) \geq \alpha$.

Algorithm 4 (*assoc_score*($\text{S}_1, T, \text{S}_2, \alpha$)) iteratively runs function *stat_independence_test*($X, T, \text{S}_2 \cup (\text{S}_1 \setminus \{X\})$), $\forall X \in \text{S}_1$. The test used by *stat_independence_test* is the G-test, to assess independence between X and T , conditional on $\text{S}_2 \cup (\text{S}_1 \setminus \{X\})$. The higher the test statistic, the higher the dependence.

In Algorithm 2 (learnMB), function *assoc_score*(s', T, MB) (line 5) is run on all subsets s' of S , including subset s , the future candidate. In particular, in Algorithm 4 (*assoc_score*), function *assoc_score*(s, T, MB) will return the maximal G-test statistic computed for some $X^* \in \text{s}$ by *stat_independence_test*($X^*, T, \text{MB} \cup (\text{s} \setminus \{X^*\})$). Let us denote test^* this test. In Algorithm 2 (learnMB), function *significant_indep*_{MB}($\text{s}, T, \text{MB}, \alpha$) (line 6) is run after function *assoc_score*(s, T, MB) (in particular) has been run. The p-value of test^* was memorized during the execution of *assoc_score*(s, T, MB) (Algorithm 4, line 4); this p-value is

then available. In learnMB, this p-value is directly used by the function $\textit{significant_indep}_{MB}(\mathbf{s}, T, \mathbf{MB}, \alpha)$, to assess independence of \mathbf{s} et T , given the current MB.