

# Automatic construction of molecular similarity networks for visual graph mining in chemical space of bioactive peptides: an unsupervised learning approach

Longendri Aguilera-Mendoza<sup>1</sup>, Yovani Marrero-Ponce<sup>2,3,4,\*</sup>, César R. García-Jacas<sup>5</sup>, Edgar Chavez<sup>1</sup>, Jesus A. Beltran<sup>1</sup>, Hugo A. Guillen-Ramirez<sup>6,7</sup>, and Carlos A. Brizuela<sup>1,\*</sup>

<sup>1</sup>Departamento de Ciencias de la Computación, Centro de Investigación Científica y de Educación Superior de Ensenada (CICESE), Baja California, 22860, Mexico

<sup>2</sup>Universidad San Francisco de Quito, Grupo de Medicina Molecular y Traslacional (MeM&T), Escuela de Medicina, Colegio de Ciencias de la Salud (COCSA), Av. Interoceánica Km 12 1/2 y Av. Florencia, 17-1200-841, Quito, Ecuador.

<sup>3</sup>Grupo GINUMED, Corporacion Universitaria Rafael Nuñez. Facultad de Salud. Programa de Medicina, Cartagena, Colombia.

<sup>4</sup>Unidad de Investigación de Diseño de Fármacos y Conectividad Molecular, Departamento de Química Física, Facultad de Farmacia, Universitat de València, Spain.

<sup>5</sup>Cátedras Conacyt - Departamento de Ciencias de la Computación, Centro de Investigación Científica y de Educación Superior de Ensenada (CICESE), Ensenada, Baja California, Mexico.

<sup>6</sup>Department of BioMedical Research (DBMR), University of Bern, Bern 3008, Switzerland

<sup>7</sup>Department of Medical Oncology, Inselspital, University Hospital and University of Bern 3010, Switzerland

\*ymarrero77@yahoo.es; cbrizuel@cicese.mx

## ABSTRACT

Algorithm SI2-1: Feature ranking and filtering.

Algorithm SI2-2: Feature subset optimization.

Algorithm SI2-3: Parallel construction of the HSP network.

## Algorithm SI2-1

---

### Algorithm 1: Feature ranking and filtering

---

```

input : A descriptor matrix  $\mathcal{D} = [x_{ij}]_{n \times m}$ , an entropy threshold  $\theta_1$ , a correlation method corr.method, and correlation-based similarity threshold  $\theta_2$ 
output : A subset  $F$  of candidate features
/* Entropy-based filtering
 $F \leftarrow \{j \mid j = 1 \dots m\}$ ;
for  $j = 1$  to  $m$  do
     $f_j.\text{entropy} \leftarrow H(f_j)$ ;
    if  $f_j.\text{entropy} < \theta_1$  then
         $| F \leftarrow F \setminus \{f_j\}$ ;
    end
end
rankedFeatures  $\leftarrow$  Sort features in  $F$  by descending order of their entropy values;
/* Correlation-based filtering
for  $j = 1$  to  $\text{sizeOf}(\text{rankedFeatures}) - 1$  do
    for  $k = j + 1$  to  $\text{sizeOf}(\text{rankedFeatures})$  do
        if corr.method = "pearson" then
             $| sim(f_j, f_k) \leftarrow |\rho(f_j, f_k)|$ ;
        else
            if corr.method = "spearman" then
                 $| sim(f_j, f_k) \leftarrow |r_s(f_j, f_k)|$ ;
            end
        end
        if  $sim(f_j, f_k) \geq \theta_2$  then
             $| F \leftarrow F \setminus \{f_k\}$ ;
        end
    end
end
end
/* Initialize the candidate set */
/* It is defined in Eq. 2 */
/* Removing irrelevant features */
/* It is defined in Eq. 5 */
/* It is defined in Eq. 6 */
/* Removing redundant features */

```

---

## Algorithm SI2-2

---

### Algorithm 2: Feature subset optimization

---

```

input : A candidate feature set  $F$ 
output : A subset  $F^*$  of optimized features
/* Second stage: subset optimization
 $F^* \leftarrow F$ ;
best.merit  $\leftarrow \Phi(F^*)$ ;
success  $\leftarrow true$ ;
while success do
    success  $\leftarrow false$ ;
    foreach  $f_j \in F^*$  do
         $F' \leftarrow F^* \setminus \{f_j\}$ ;
        merit  $\leftarrow \Phi(F')$ ;
        if merit > best.merit then
            best.merit  $\leftarrow merit$ ;
            best.subset  $\leftarrow F'$ ;
            success  $\leftarrow true$ ;
        end
    end
    if success then
         $| F^* \leftarrow best.\text{subset}$ ;
    end
end
/* Candidate feature set */
/* It is defined in Eq. 7 */

```

---

## Algorithm SI2-3

**Algorithm 3:** Parallel construction of the HSP network

---

```

input : A descriptor matrix  $\mathcal{D} = [x_{ij}]_{n \times m}$ , and distance function  $d$ 
output : A weighted graph  $G' = (V, E', w)$  with a weight  $w : E' \rightarrow [0, 1]$ 
 $V \leftarrow \{i \mid i = 1 \dots n\};$ 
 $E' \leftarrow \emptyset;$ 
 $maxDist \leftarrow 0;$ 
for  $u = 1$  to  $n$  do in parallel
    candidates  $\leftarrow$  Array[0.. $n - 1$ ];
    cursor  $\leftarrow 0$ ;
    foreach  $v \in V$  do
        if  $v \neq u$  then
            dist  $\leftarrow d(u, v)$ ;
            candidates[cursor]  $\leftarrow$  Candidate( $v, dist$ );
            cursor  $\leftarrow cursor + 1$ ;
        end
    end
    Sort candidates by ascending order of their distances to node  $u$ ;
    largeDist  $\leftarrow$  candidates[ $n - 1$ ].distance;
    cursor  $\leftarrow 0$ ;
    while cursor <  $n - 1$  do
        if candidates[cursor]  $\neq$  null then
            v  $\leftarrow$  candidates[cursor].node;
            dist  $\leftarrow$  candidates[cursor].distance;
            writeLock();
             $E' \leftarrow E' \cup \{(u, v)\}$ ;
            writeUnlock();
            /* Ignoring candidates in the forbidden area for node  $u$ 
            for  $k = cursor + 1$  to  $n - 1$  do
                if candidates[k]  $\neq$  null then
                    dist  $\leftarrow d(v, candidates[k].node)$ ;
                    if dist < candidates[k].distance then
                        candidates[k]  $\leftarrow$  null;
                    end
                end
            end
            cursor  $\leftarrow cursor + 1$ ;
        end
        writeLock();
        if largeDist > maxDist then
            maxDist  $\leftarrow$  largeDist;
        end
        writeUnlock();
    end
    foreach  $(u, v) \in E'$  do
         $w(u, v) \leftarrow sim(u, v) \leftarrow 1 - \frac{d(u, v)}{maxDist}$ ;
    end
    /* The similarity defined in Eq. 8 */

```

---