

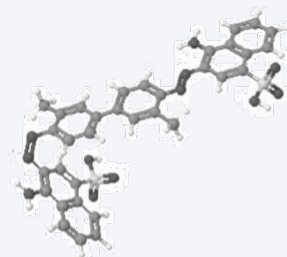
# Task 5: Molecules

AIDS Antiviral Screen Database of Active Compounds

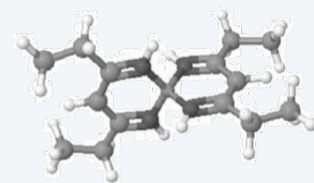
250 training, 250 validation molecules

Two classes active 'a' and inactive 'i'

Annotated in `train.txt` and `valid.txt`



(a) Active



(b) Inactive

## Task

Classify the molecules of the validation set using KNN

Distance: approximate Graph Edit Distance (GED)

# Input: Graph xml (gxl files)

## XMLs with a lot of information

→ we only need nodes labeled with their chemical `symbol` and the unlabeled, undirected edges:

```
<node id="_1"><attr name="symbol"><string>C</string></attr>
```

...

```
<edge from="_1" to="_2">
```

...

## Further info about the database (AIDS):

See `riesen08graphdb.pdf` → section 2.8 (there are also images of molecules)

Hint: there are python libraries to parse XML files

# Task 5: Molecules

Compute approximate GED between pairs of molecules with  
**bipartite graph matching**

(lecture 10, slide 21)

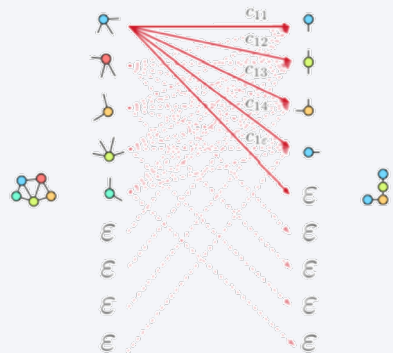
Build cost matrix (*Dirac*)

$$C = \left[ \begin{array}{cccc|cccc} c_{11} & c_{12} & \dots & c_{1m} & c_{1e} & \infty & \dots & \infty \\ c_{21} & c_{22} & \dots & c_{2m} & \infty & c_{2e} & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \infty \\ c_{n1} & c_{n2} & \dots & c_{nm} & \infty & \dots & \infty & c_{ne} \\ \hline c_{e1} & \infty & \dots & \infty & 0 & 0 & \dots & 0 \\ \infty & c_{e2} & \dots & \vdots & 0 & 0 & \dots & \vdots \\ \vdots & \vdots & \ddots & \infty & \vdots & \vdots & \ddots & 0 \\ \infty & \dots & \infty & c_{em} & 0 & \dots & 0 & 0 \end{array} \right]$$

Hungarian Algorithm

To find optimal assignment

Derive Edit Path costs from the result  
(distance for classification)



KNN for classification (optimize for  $\mathbb{K}$ )

# Task 5: Molecules

## Recommendation

Use *Dirac* cost function for GED (optimize  $C_n$  and  $C_e$ )  
(lecture 9, slide 36)

Node substitution:  $2 * C_n$  if symbols  $\neq$ , 0 otherwise

Node deletion/insertion:  $C_n$

Edge deletion/insertion:  $C_e$

Use an existing framework for the Hungarian algorithm