

Quimioinformática en Phyton

Algoritmos de clustering, métricas de validación y aplicaciones, y métricas de enriquecimiento, campañas de cribado retrospectivo y generación de señuelos

Alan Talevi

Laboratorio de Investigación y Desarrollo de Bioactivos (LIDeB)

Universidad Nacional de La Plata – CONICET

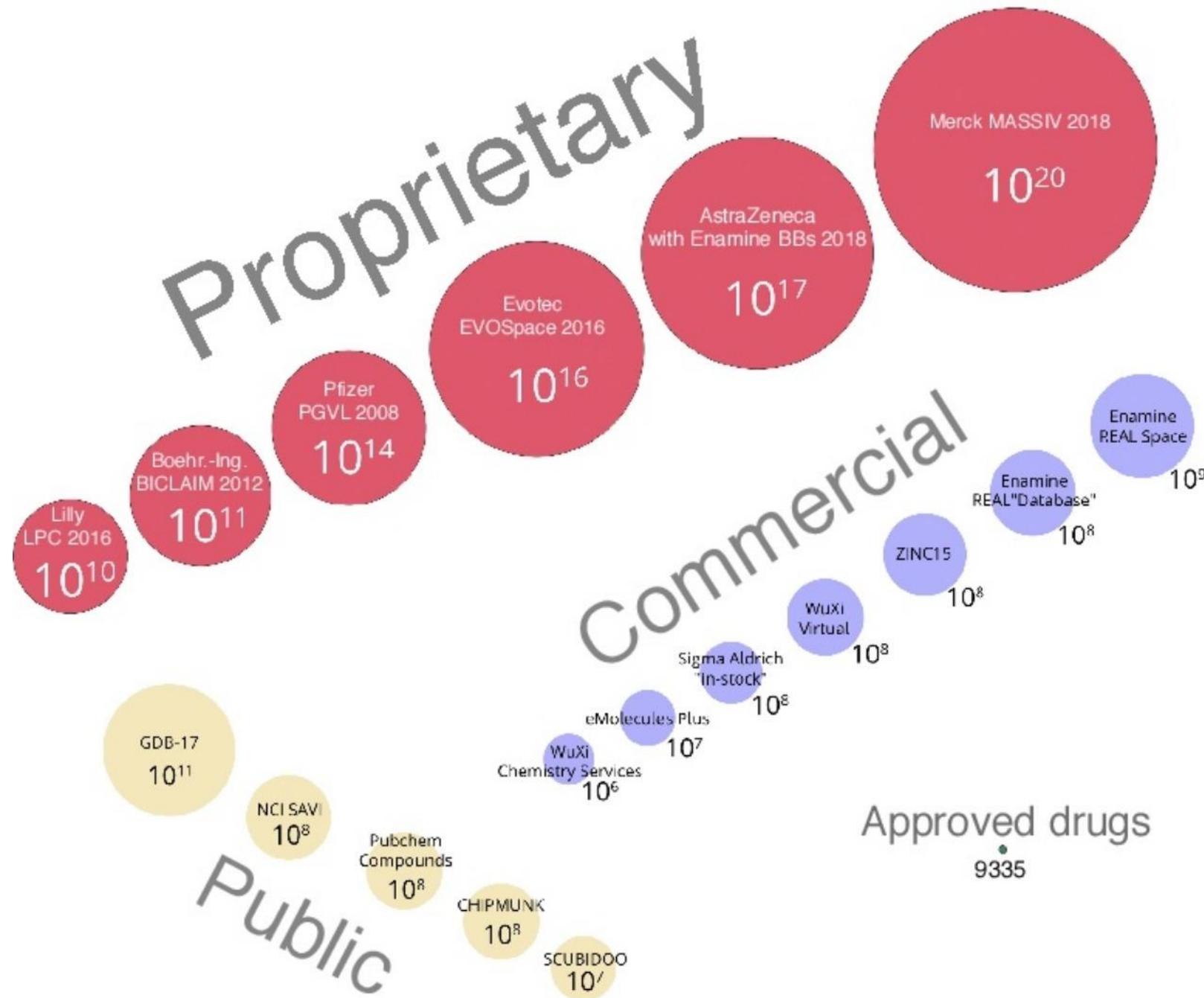
## Hoja de ruta

Contexto. Cribado virtual como estrategia eficiente para explorar el espacio químico.  
Andamiaje activo y el hit-to-lead.

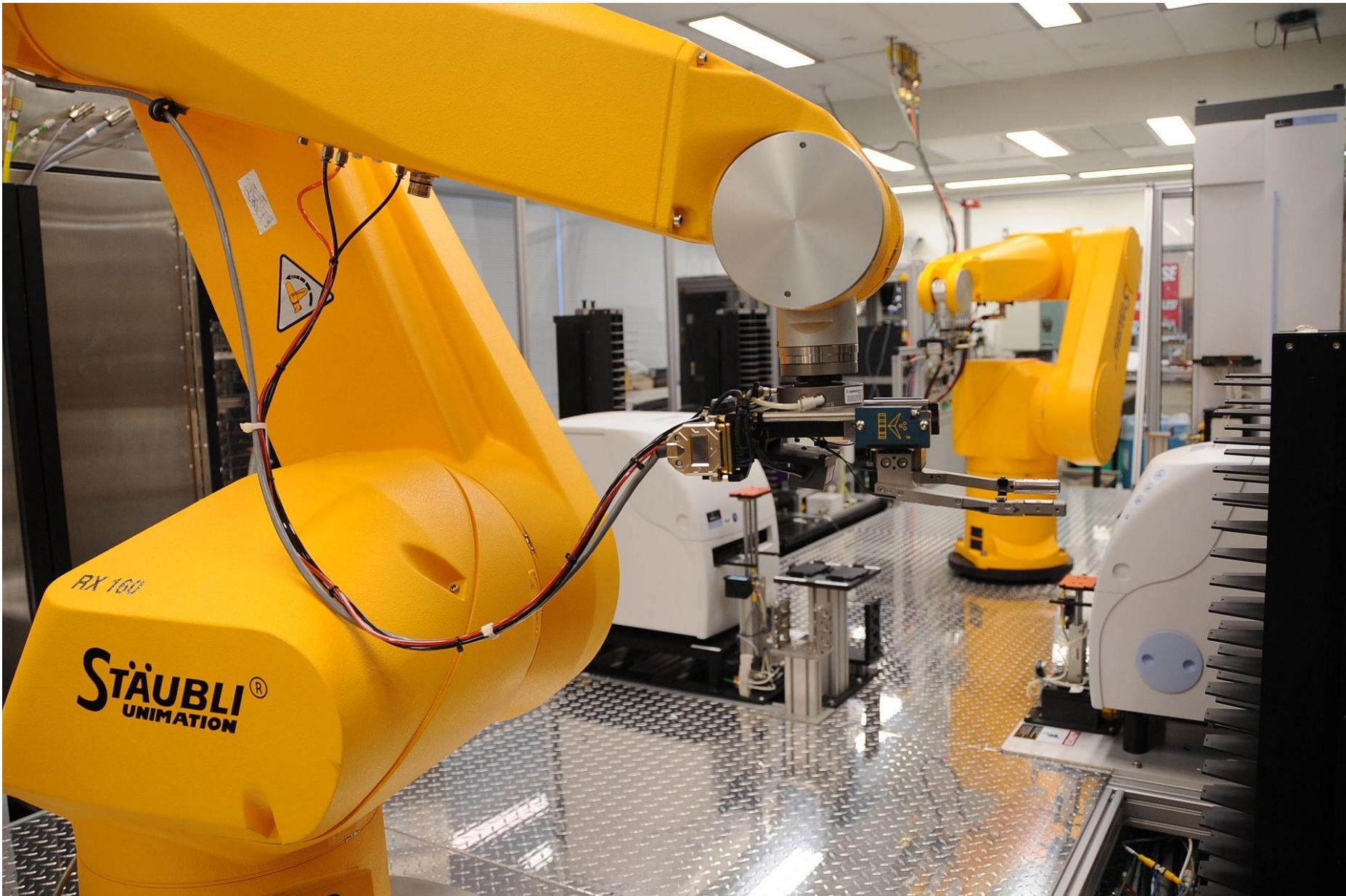
El problema del desbalance de clases en el cribado virtual. Importancia de las campañas de cribado retrospectivo para evaluar el desempeño. Métricas de enriquecimiento.

Generación de señuelos.

Clustering de pequeñas moléculas.



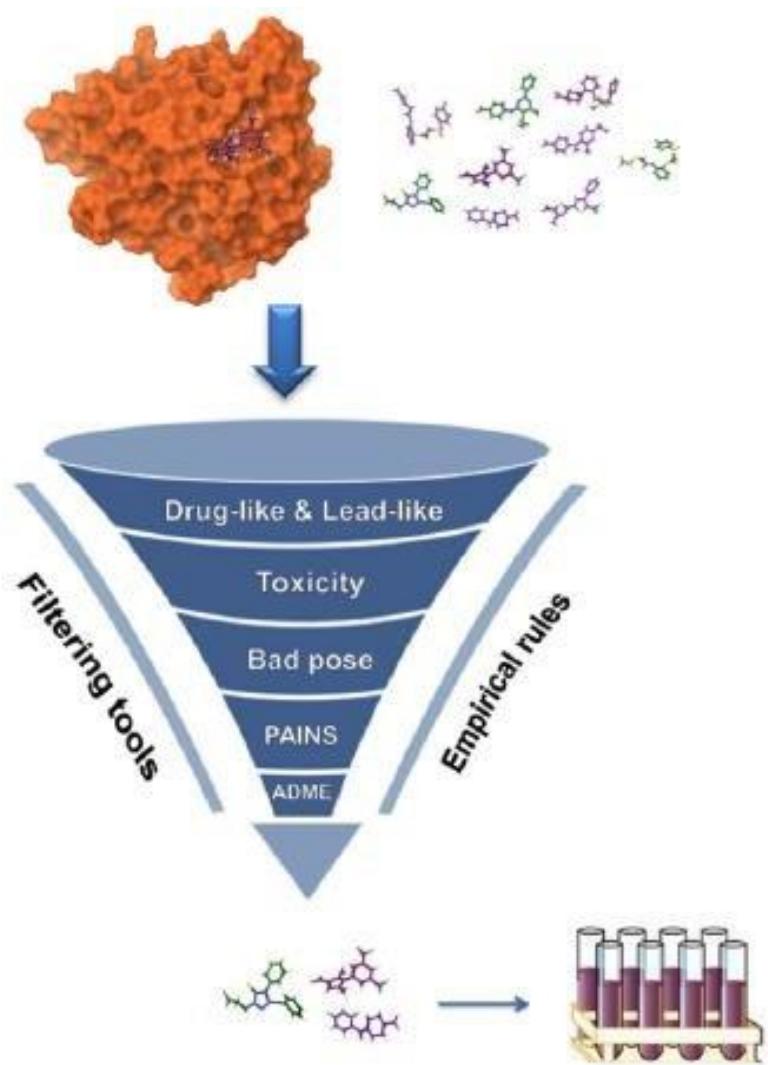




## virtual screening

utilizar modelo(s) computacional(es)/algoritmo(s) para identificar en [grandes] quimiotecas digitales, compuestos que posean una o más propiedades [biológicas] de interés (y, ocasionalmente, una diversidad de propiedades farmacéuticamente relevantes).

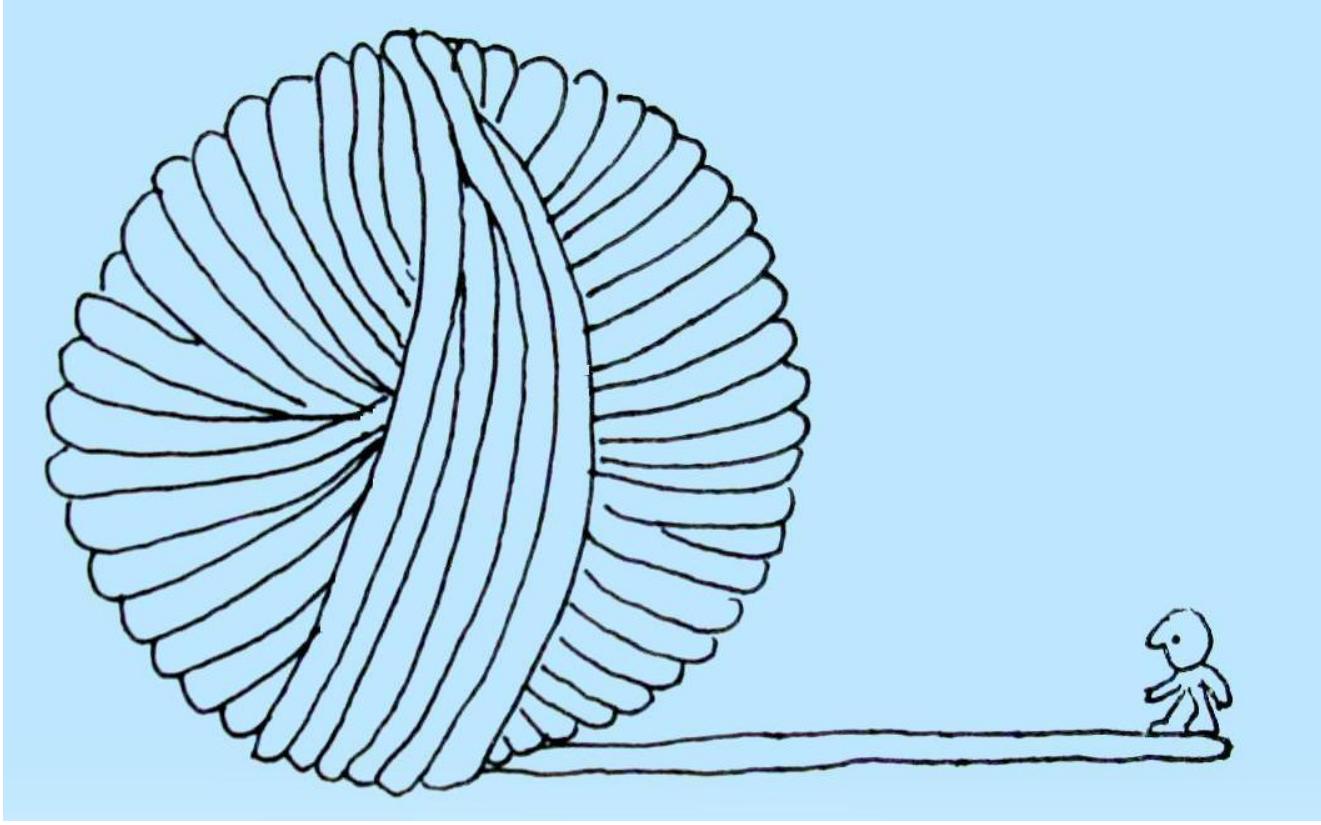
Por lo general, el hit in silico confirmado se someterá a una campaña de optimización molecular (hit-to-lead).



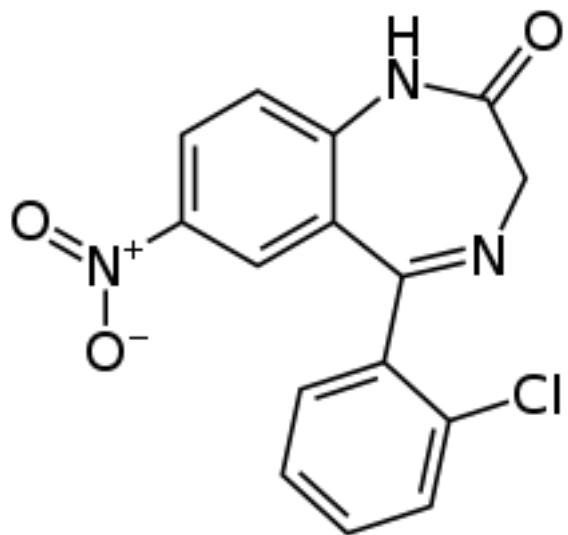
algunas ventajas del screening in silico



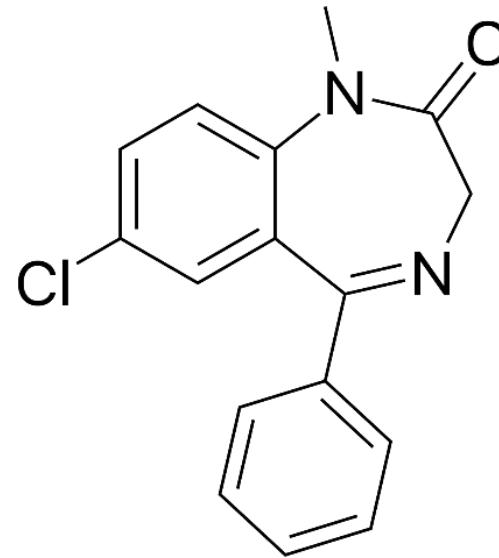
¿Cómo se descubre un nuevo fármaco?



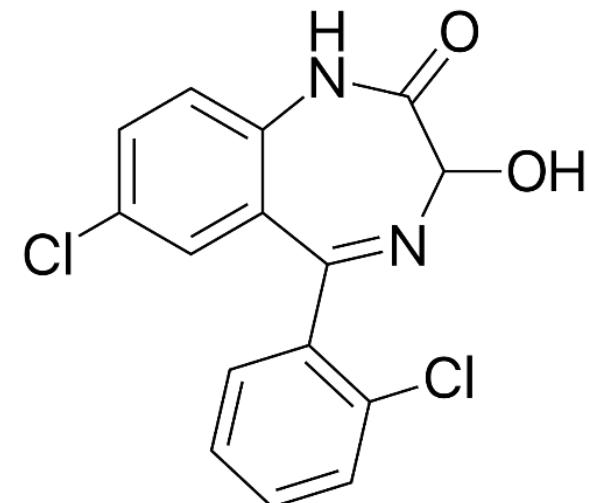
Andamiaje activo ("active scaffold"):  
núcleo químico con determinada actividad biológica



clonazepam

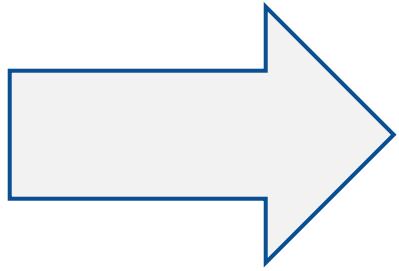


diazepam



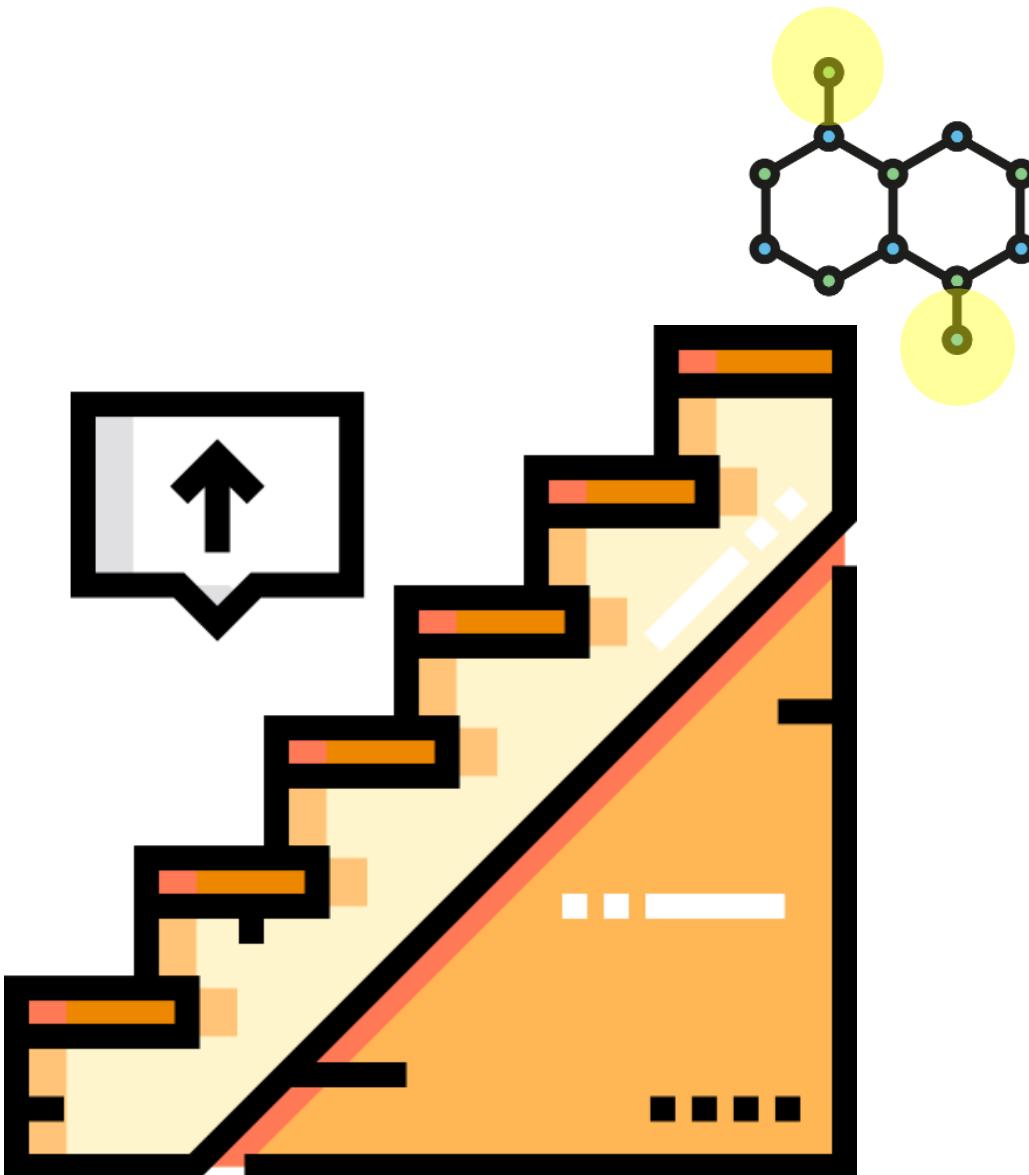
lorazepam

HIT



LÍDER O CABEZA  
DE SERIE

# Impacto de la optimización de hits y líderes



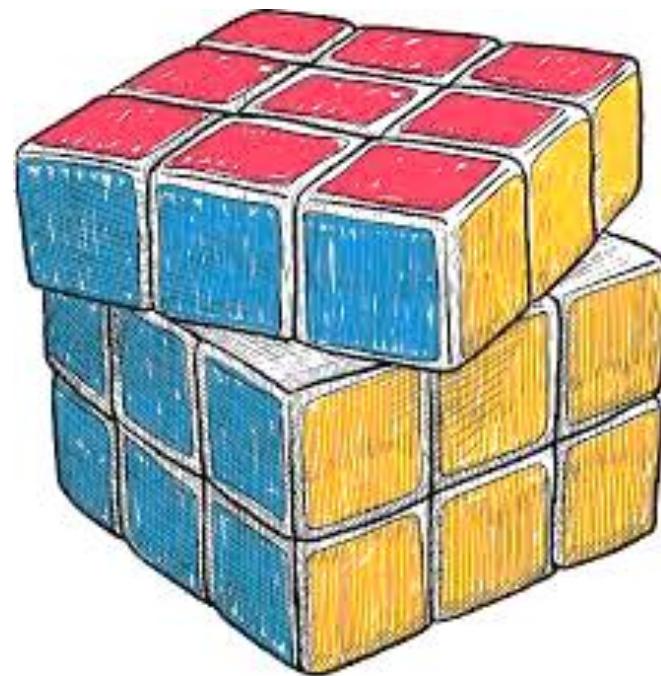
Incrementar la potencia / afinidad hacia el blanco farmacológico

Eliminación de efectos adversos

Mejorar biodisponibilidad

Aumentar la estabilidad química

Aumentar la novedad estructural



# El problema del desbalance de clases

$$P(A) = (Se \cdot Ya) / [Se \cdot Ya + (1-Sp) (1-Ya)]$$

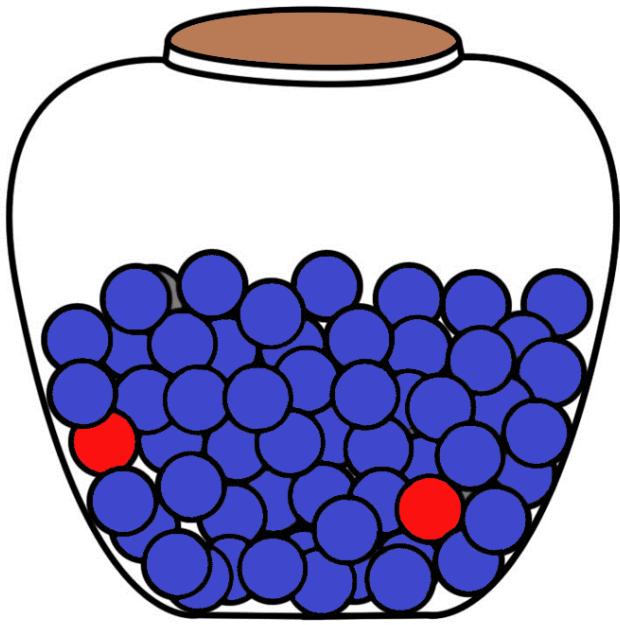
$P(A)$ : probabilidad de que un hit in silico confirmará experimentalmente la actividad predicha

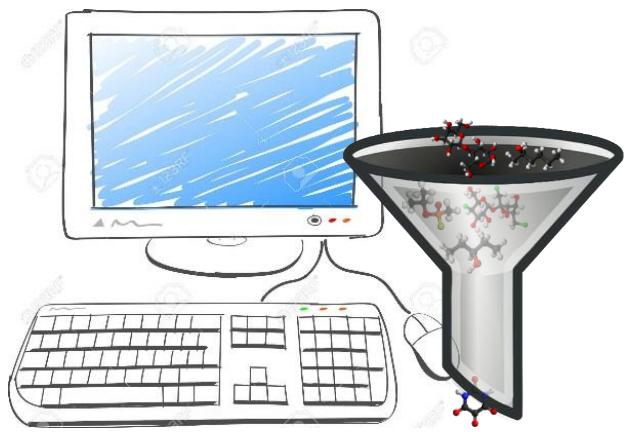
$Se$ : Sensibilidad, tasa de verdaderos positivos

$Sp$ : Especificidad, tasa de verdaderos negativos

$Ya$  (o  $Ra$ ): proporción de compuestos activos en la quimioteca sometida a screening (**desconocida a priori**)

Ya empírico 0,005 - 0,01  
(para quimiotecas no enfocadas)





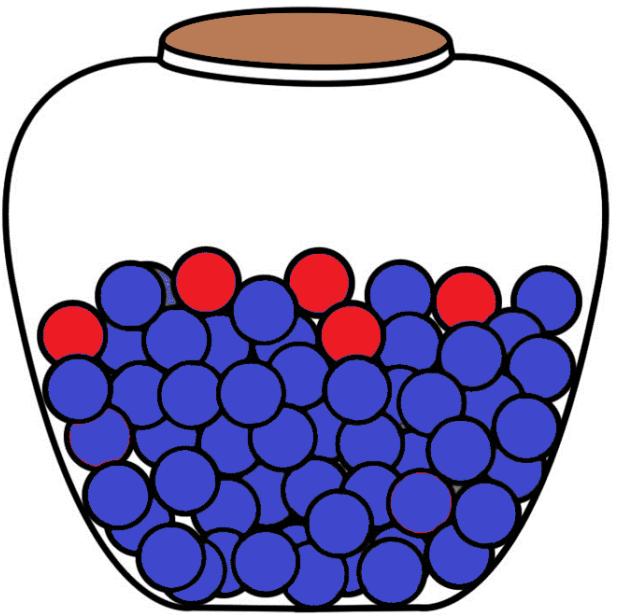
INACTIVO
INACTIVO
INACTIVO
ACTIVO
INACTIVO
ACTIVO
...

cribado  
virtual



ACTIVO
INACTIVO
ACTIVO
INACTIVO
ACTIVO
INACTIVO
ACTIVO
ACTIVO
INACTIVO
ACTIVO
...

si el único criterio de selección fuera el score del modelo, enviaríamos a confirmación experimental las primeras posiciones del ranking ordenado



La capacidad de enriquecimiento de un protocolo de cribado virtual se estima a través de experimentos retrospectivos, en los que se siembra un número relativamente bajo (pero diverso) de compuestos activos conocidos entre un número elevado de señuelos o *decoys* (compuestos inactivos verificados o presuntos/putativos). En el experimento retrospectivo se conocen la identidad y proporción ( $Y_a / R_a$ ) de los compuestos activos



Métricas de enriquecimiento utilizadas  
en experimentos de cribado retrospectivo

## Métricas de enriquecimiento

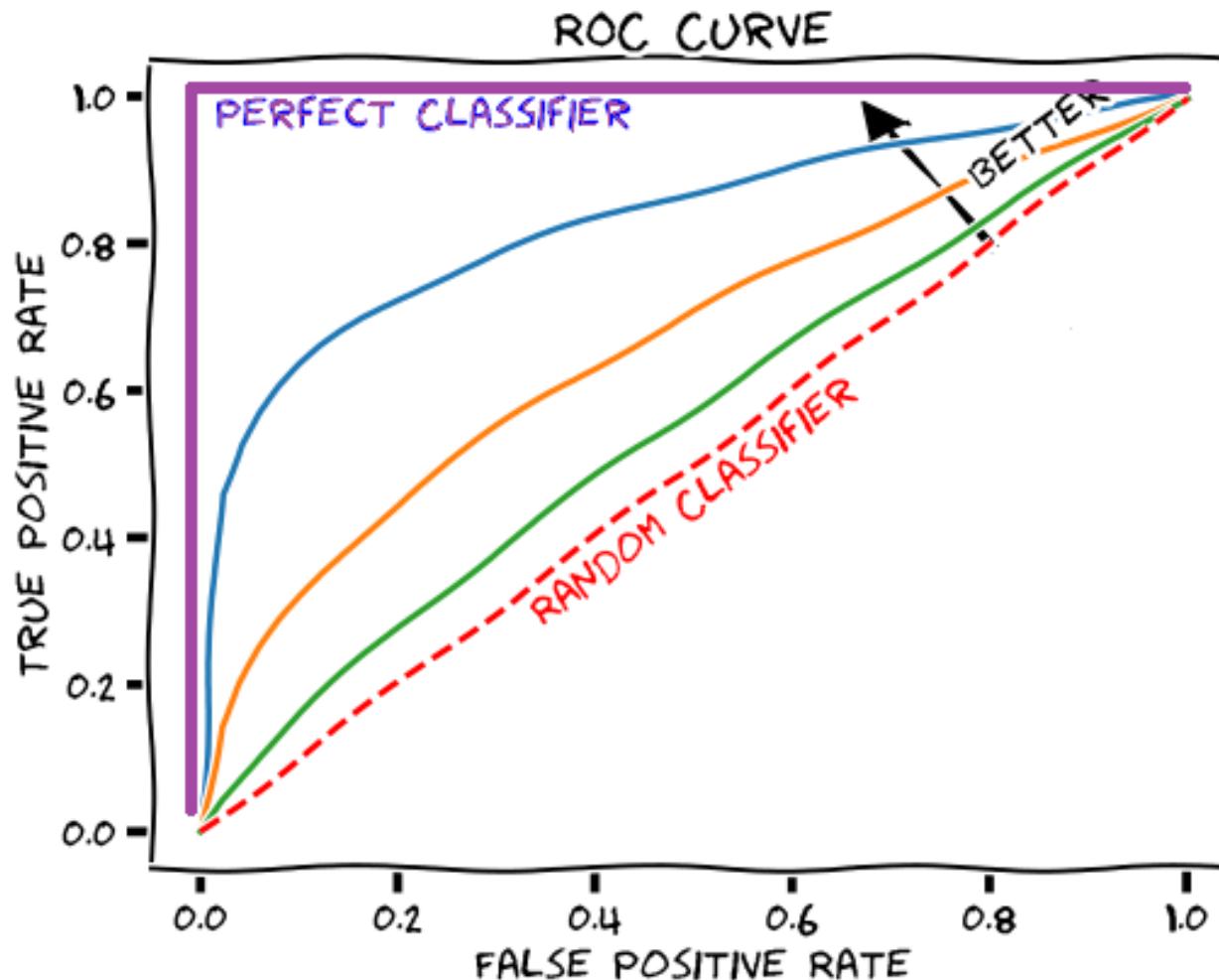
Permiten evaluar el desempeño de un modelo o protocolo de cribado virtual y comparar cuantitativamente / estadísticamente modelos o protocolos.

Sheridan, R. P.; Singh, S. B.; Fluder, E. M.; Kearsley, S. K. Protocols for bridging the peptide to nonpeptide gap in topological similarity searches. *J. Chem. Inf. Comput. Sci.* 2001

Triballeau N, Acher F, Brabet I, Pin JP, Bertrand HO. Virtual screening workflow development guided by the "receiver operating characteristic" curve approach. Application to high-throughput docking on metabotropic glutamate receptor subtype 4. *J Med Chem.* 2005

Truchon JF, Bayly CI. Evaluating virtual screening methods: good and bad metrics for the "early recognition" problem. *J Chem Inf Model.* 2007

## Área bajo la Receiver Operating Characteristic (ROC) curve - AUROC

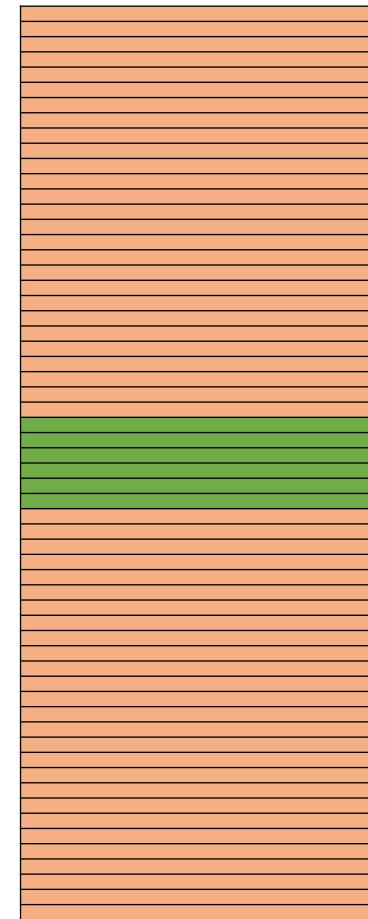
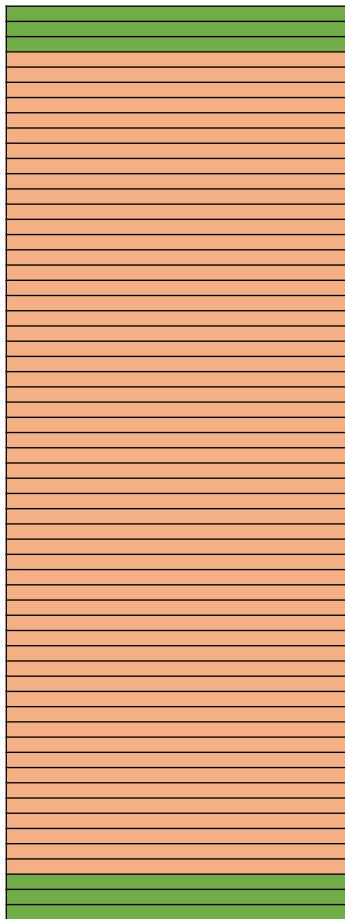
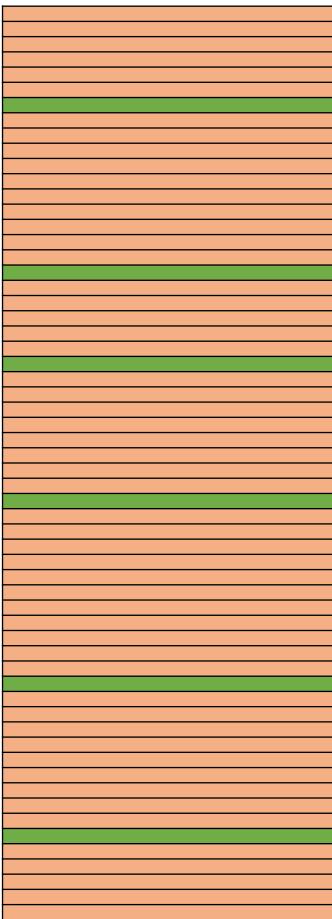


## Área bajo la curva ROC (AUROC, ROC): limitaciones

AUROC / ROC refleja la posición promedio de los verdaderos positivos en el ranking ordenado. Todas las posiciones del ranking contribuyen al valor de la métrica de igual manera.

Efecto de saturación.

## AUROC: limitaciones



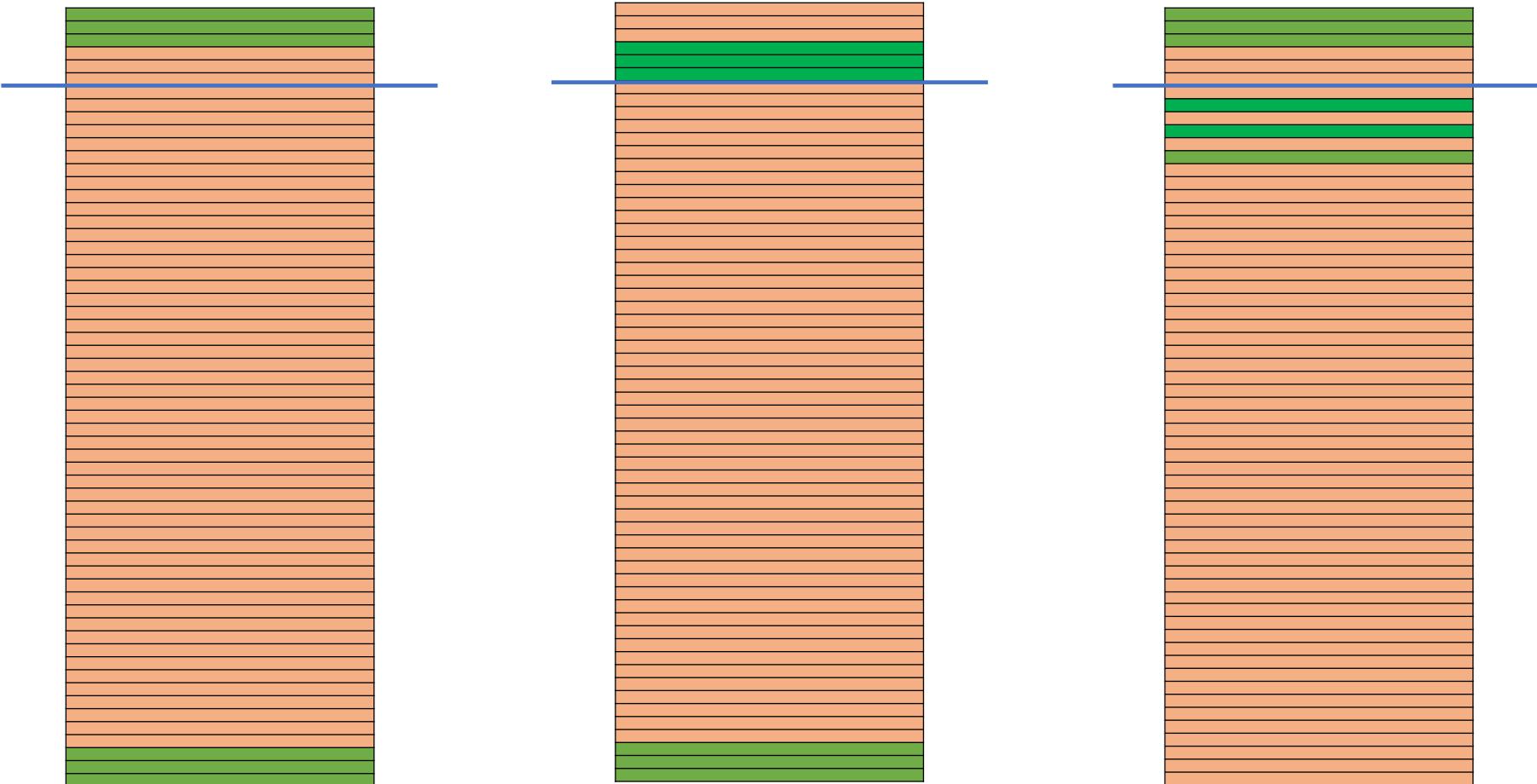
$$N = 60; n = 6; Ra = 6/60 = 0,1$$

Los tres rankings ordenados que se muestran resultan en una AUROC  $\approx 0,5$

## Métricas enfocadas en el reconocimiento temprano: Factor de enriquecimiento (EF)

$$EF_{x\%} = \frac{\text{VP en el } x\% \text{ de la quimioteca mejor rankeado por el modelo}}{\text{VP en el } x\% \text{ mejor rankeado de la quimioteca ordenada al azar}}$$

## ER: limitaciones



$$N = 60; n = 6; Ra = 6/60 = 0,1$$

¡En los tres casos el  $ER_{10\%}$  es idéntico  $\approx 5!$

## ER: limitaciones

No distingue diferencias entre distintos ordenamientos de los activos dentro del  $x\%$  de la quimioteca en el que se pone el foco.

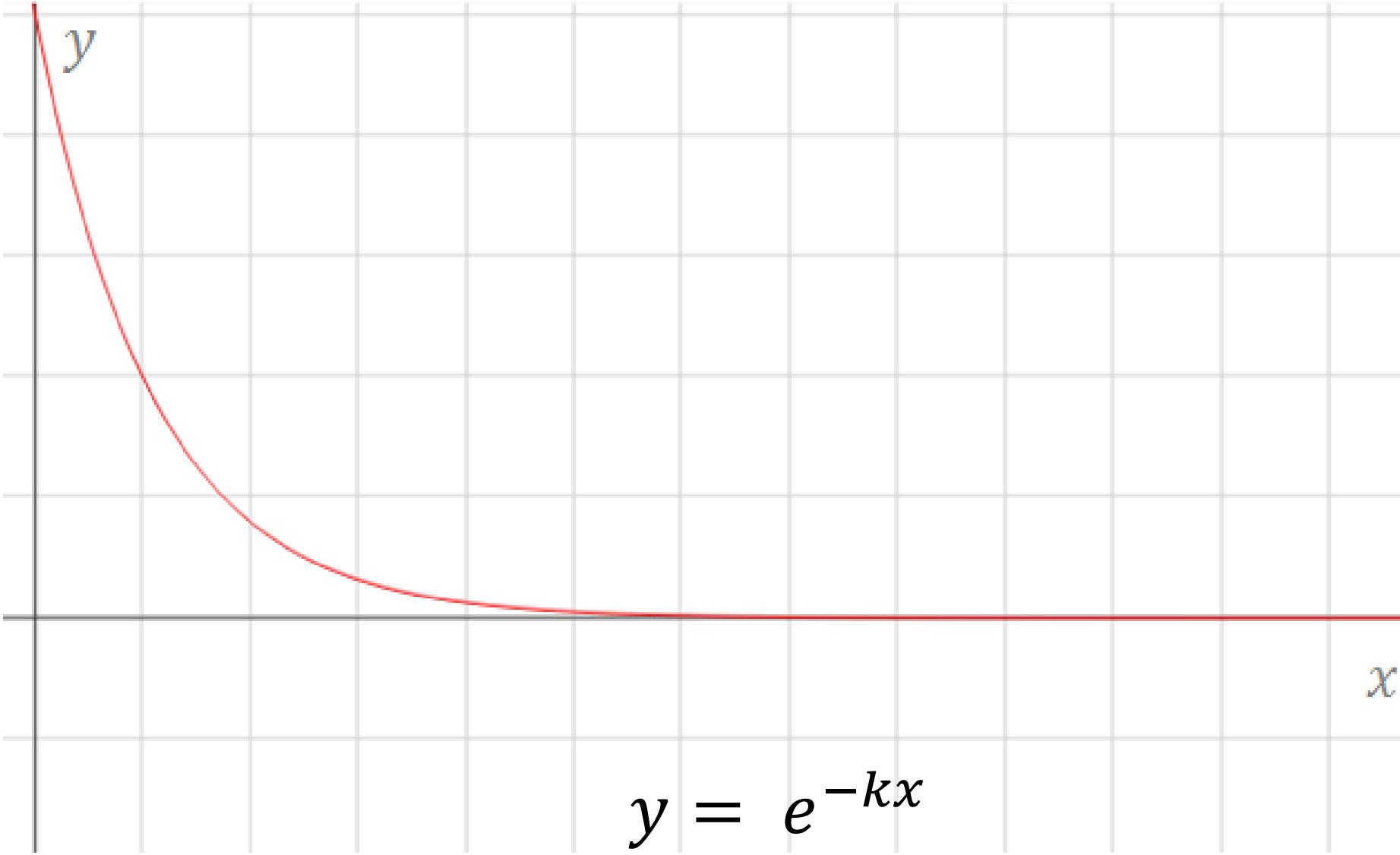
No distingue diferencias entre distintos ordenamientos de los activos en el  $100-x\%$  de la quimioteca en el que no se hace foco. Solo considera el  $x\%$  mejor rankeado.

La cota superior de  $EF_x\%$  depende de  $x$ , de  $n$  y de  $N$ . Es  $1/x$  si  $x \geq n/N$  y  $N/n$  si  $x < n/N$ .

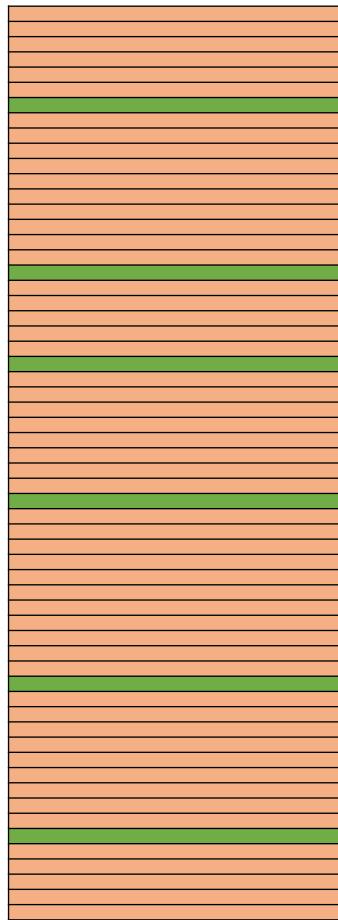
## Métricas enfocadas en el reconocimiento temprano: Robust Initial Enhancement (RIE)

$$RIE = \frac{\sum_1^n e^{-\alpha r_i/N}}{\langle \sum_1^n e^{-\alpha r_i/N} \rangle_r}$$

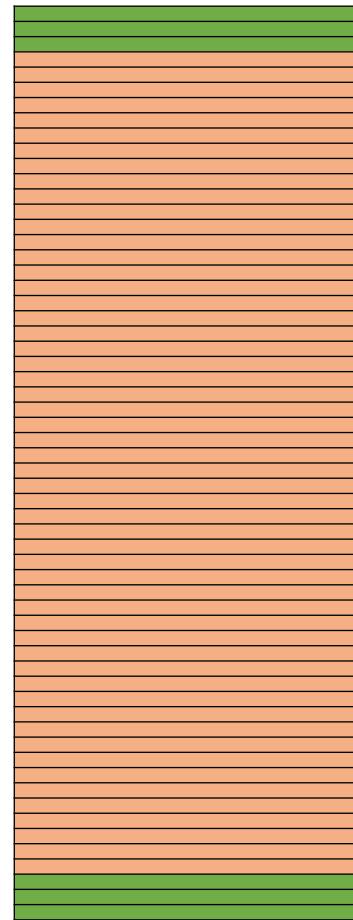
(ri = posición en el ranking del i-ésimo compuesto activo)



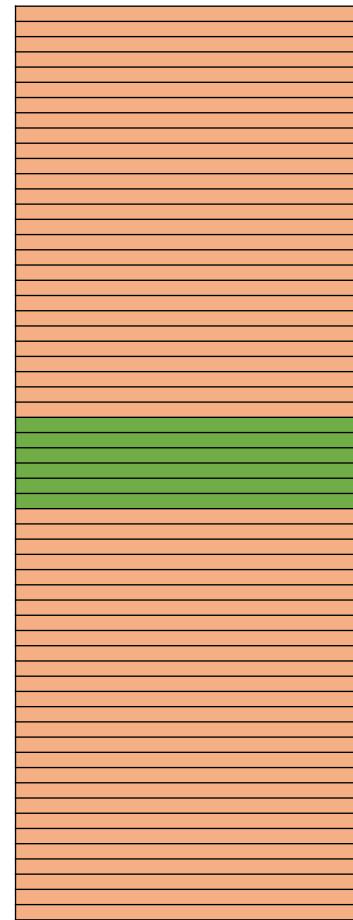
a)



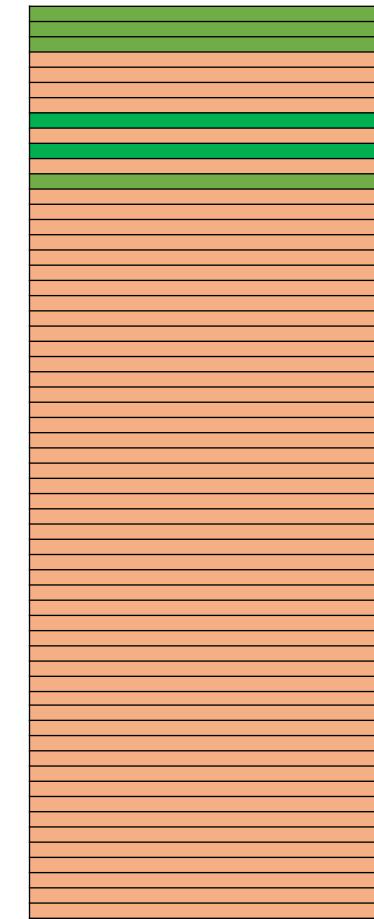
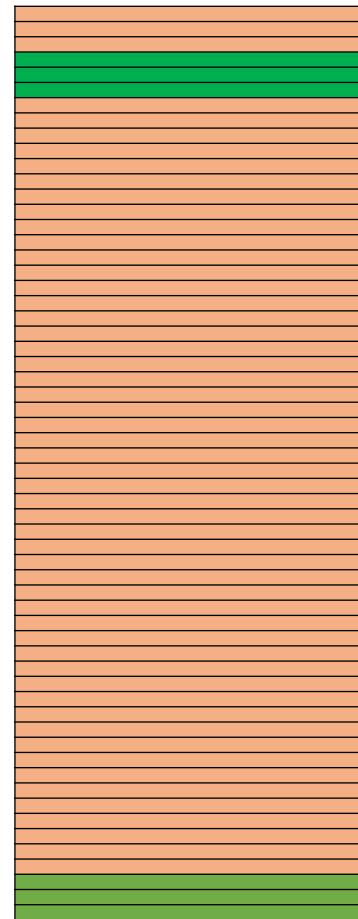
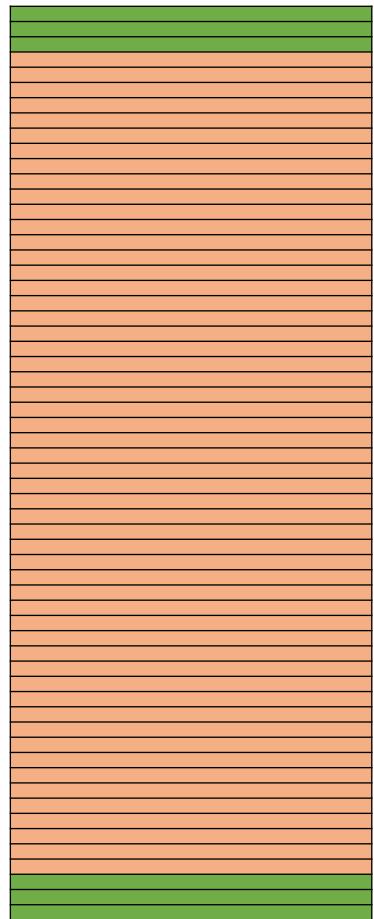
b)



c)



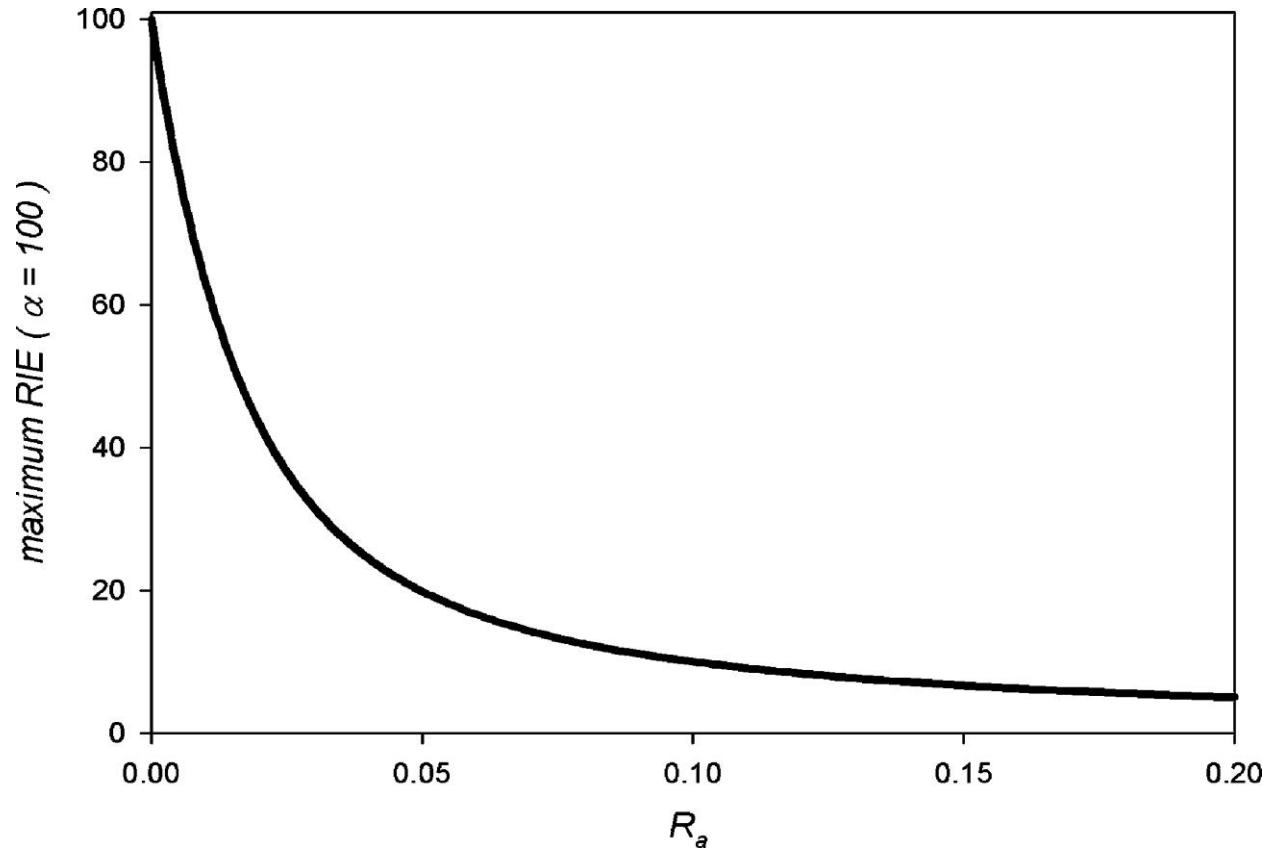
Estos tres ordenamientos posibles tendrían valores de RIE **muy** distintos.  $RIE_b > RIE_a > RIE_c$



Estos tres ordenamientos posibles también tendrían valores de RIE distintos.

## RIE: limitaciones

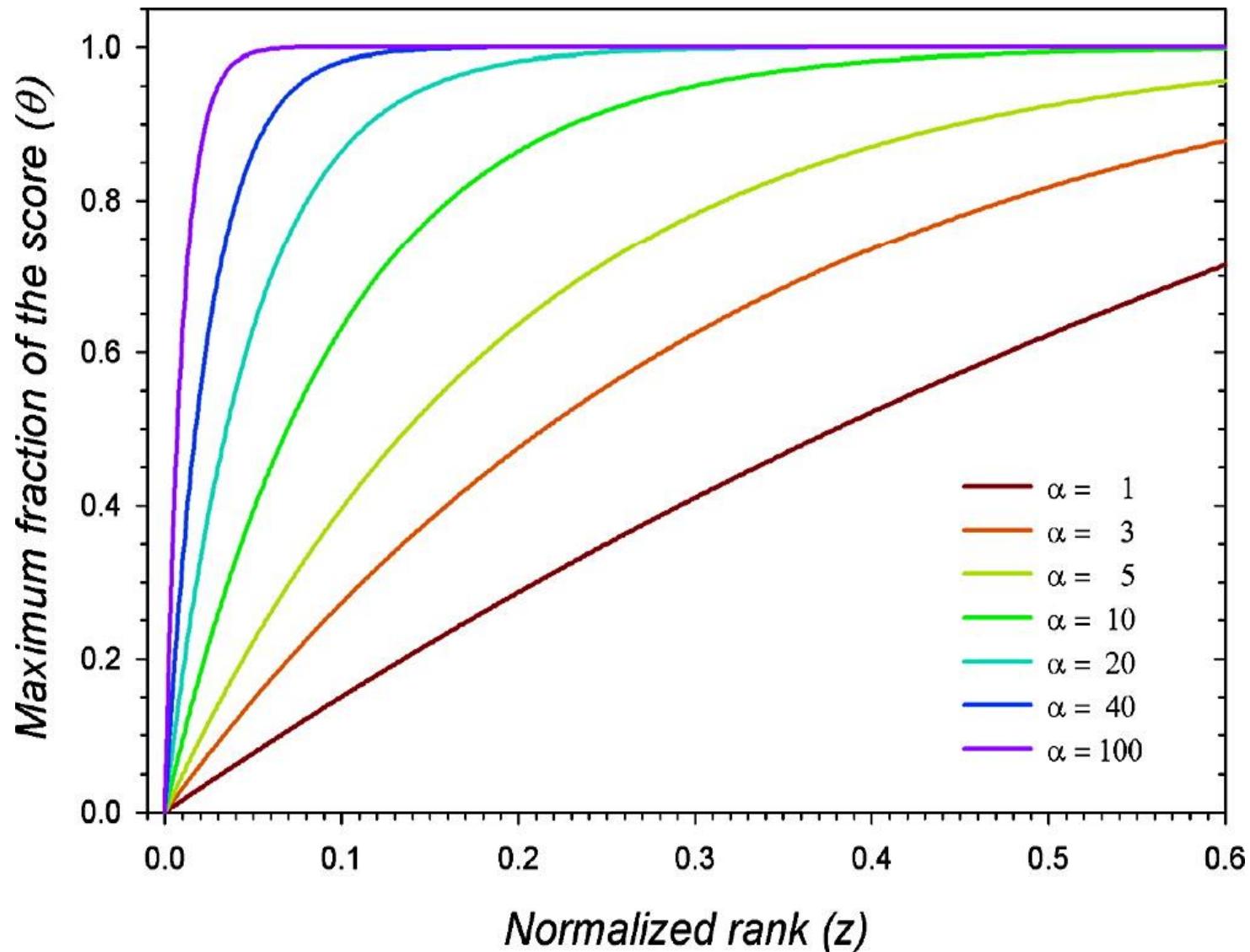
Los valores máximo y mínimo de RIE dependen de  $\alpha$ , n y N.

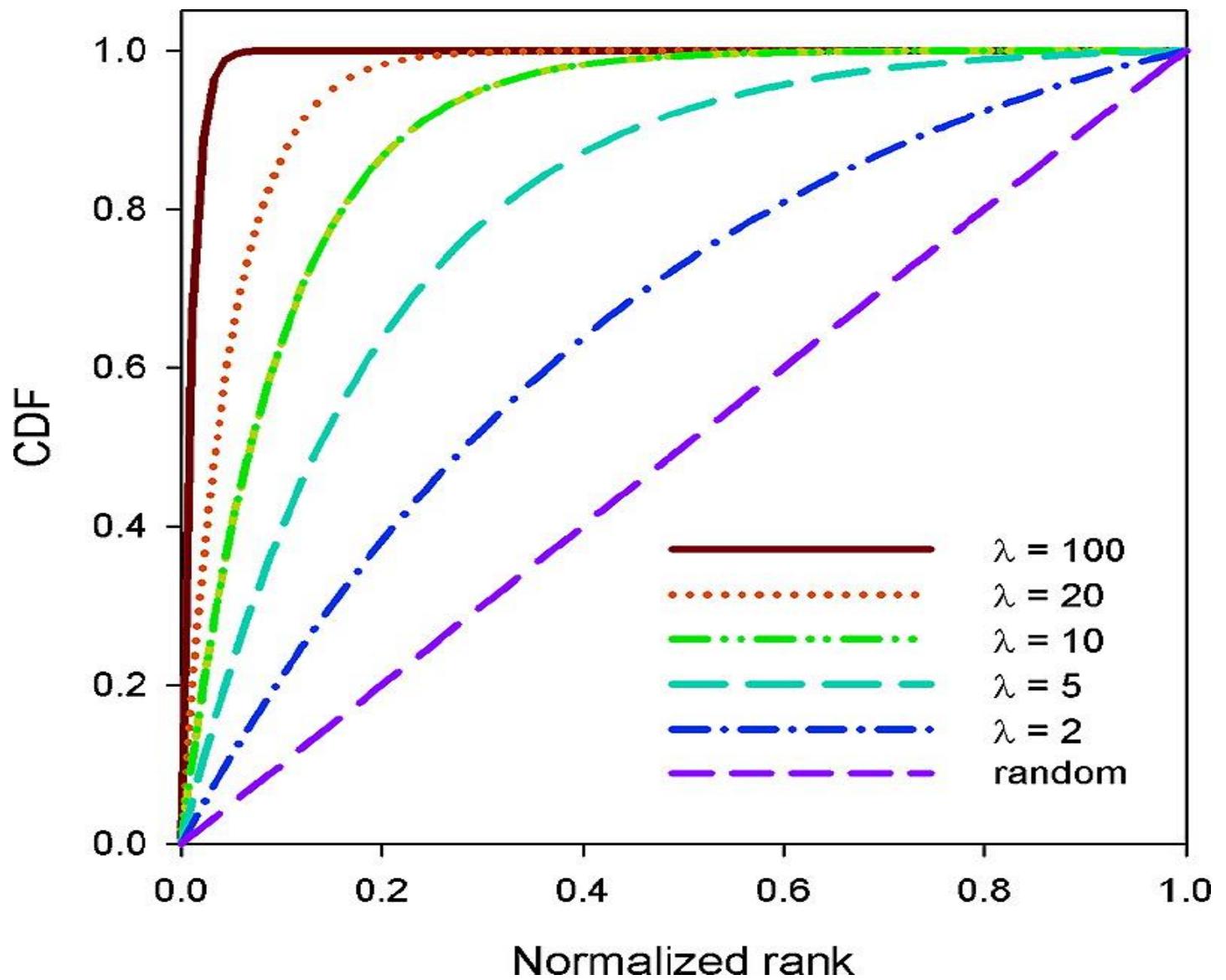


## Métricas enfocadas en el reconocimiento temprano: Boltzmann-Enhanced Discrimination of ROC (BEDROC)

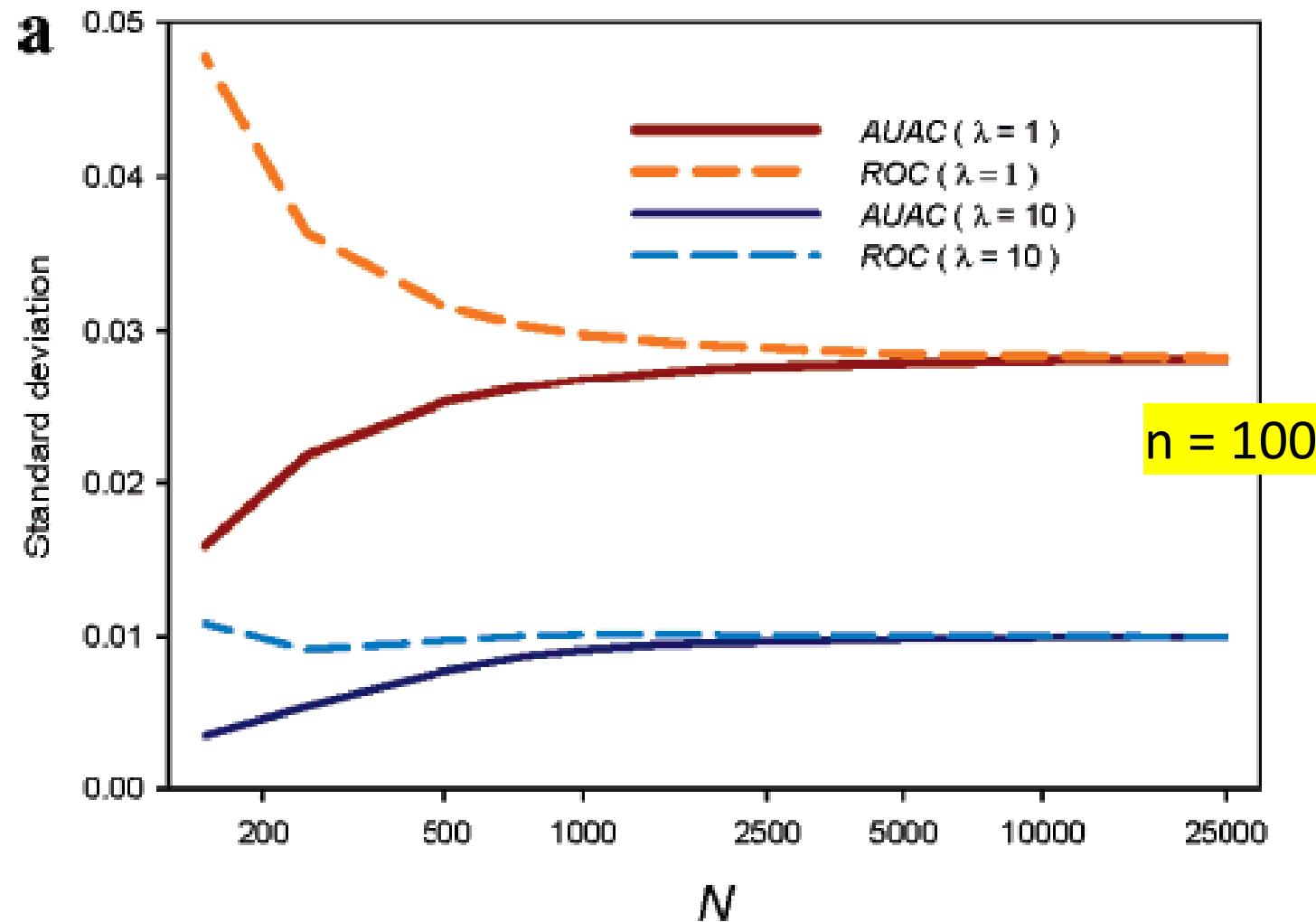
$$BEDROC = \frac{RIE - RIE_{min}}{RIE_{max} - RIE_{min}}$$

Relación entre el valor  $\alpha$  y la fracción mejor rankeada de la lista ordenada que contribuye en mayor medida al valor de la métrica BEDROC. Por ejemplo, si  $\alpha = 100$ , el 1% mejor rankeado contribuye aprox. con el

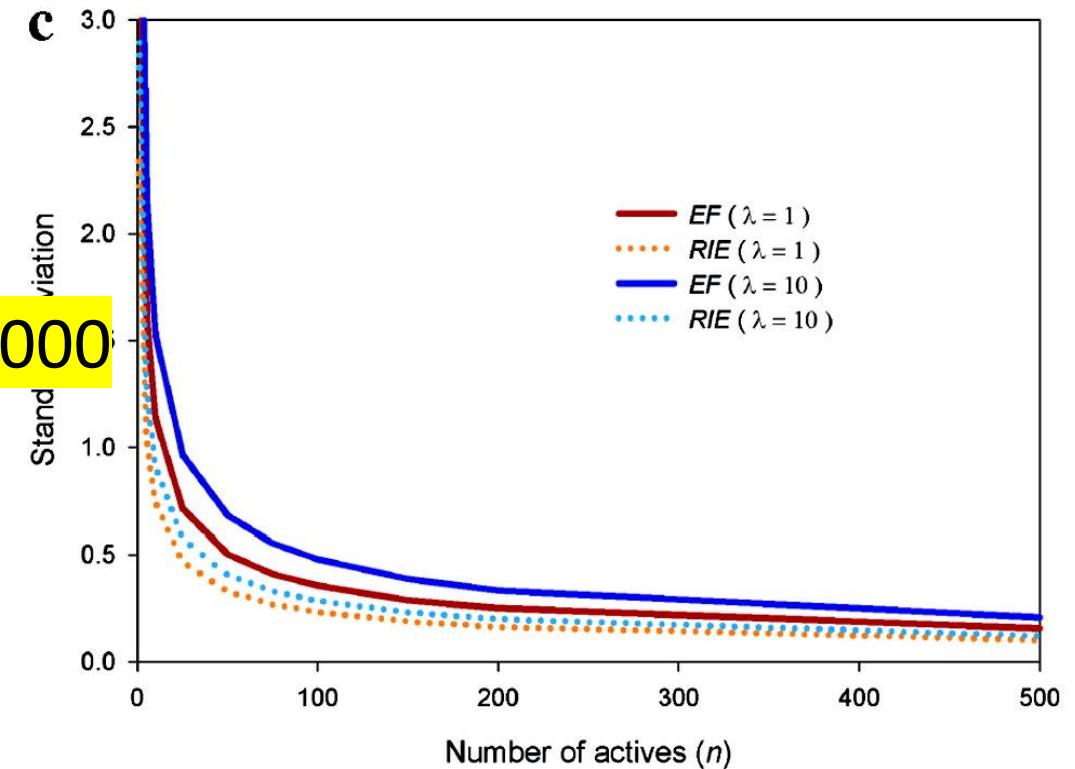
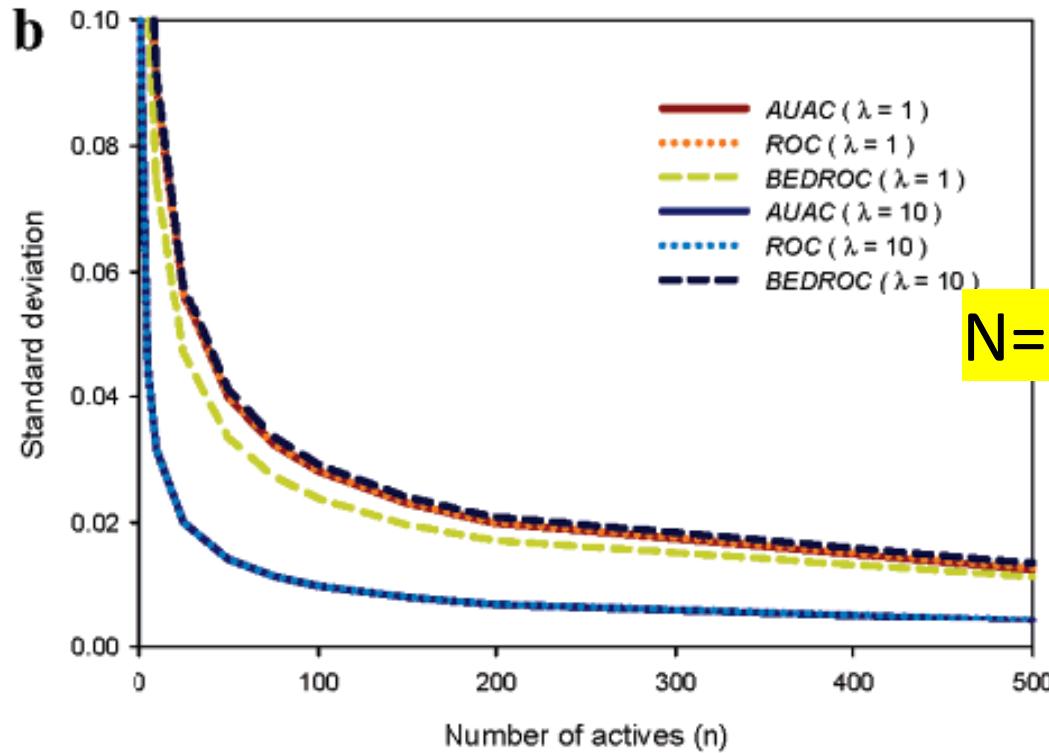




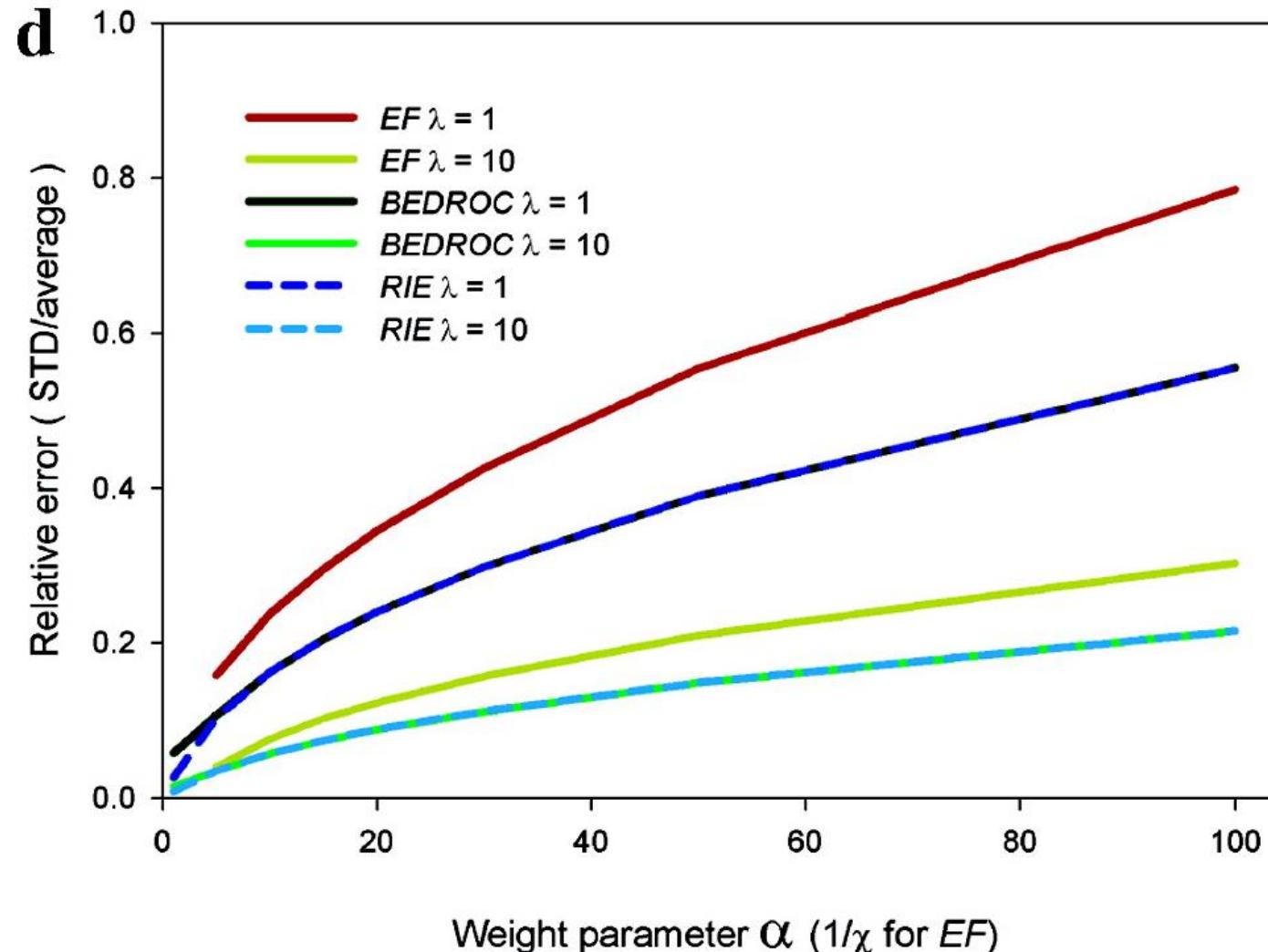
# Desviación estándar de ROC versus N, con n constante



# Desvío estándar de métricas de enriquecimiento versus n, N constante

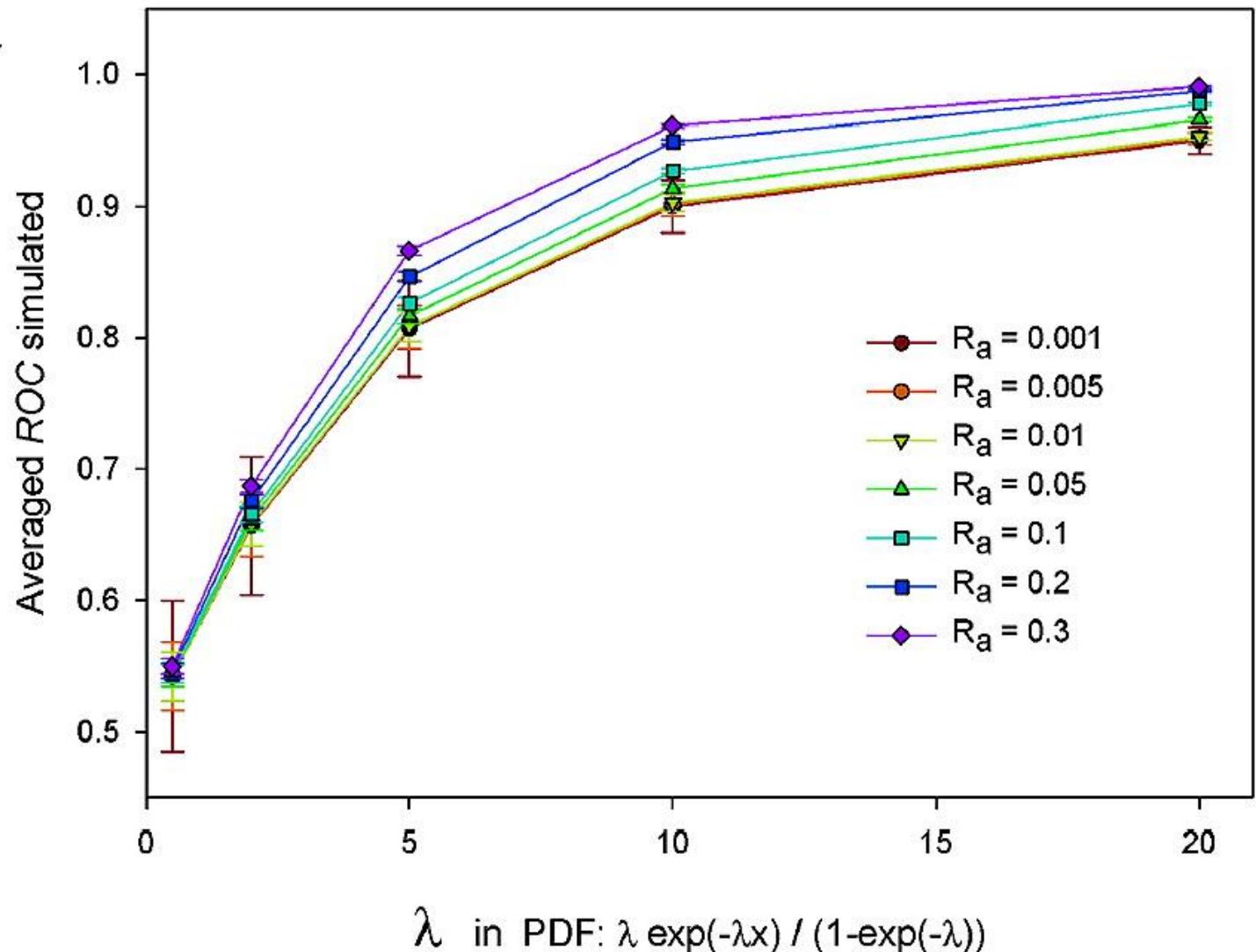


# Error relativo de métricas de enriquecimiento temprano versus $\alpha$ (RIE, BEDROC) o $1/x$ (EF). $\lambda$ variable



# Efecto de saturación

a



Para prevenir efecto de saturación:

ROC: trabajar con  $R_a << 1$

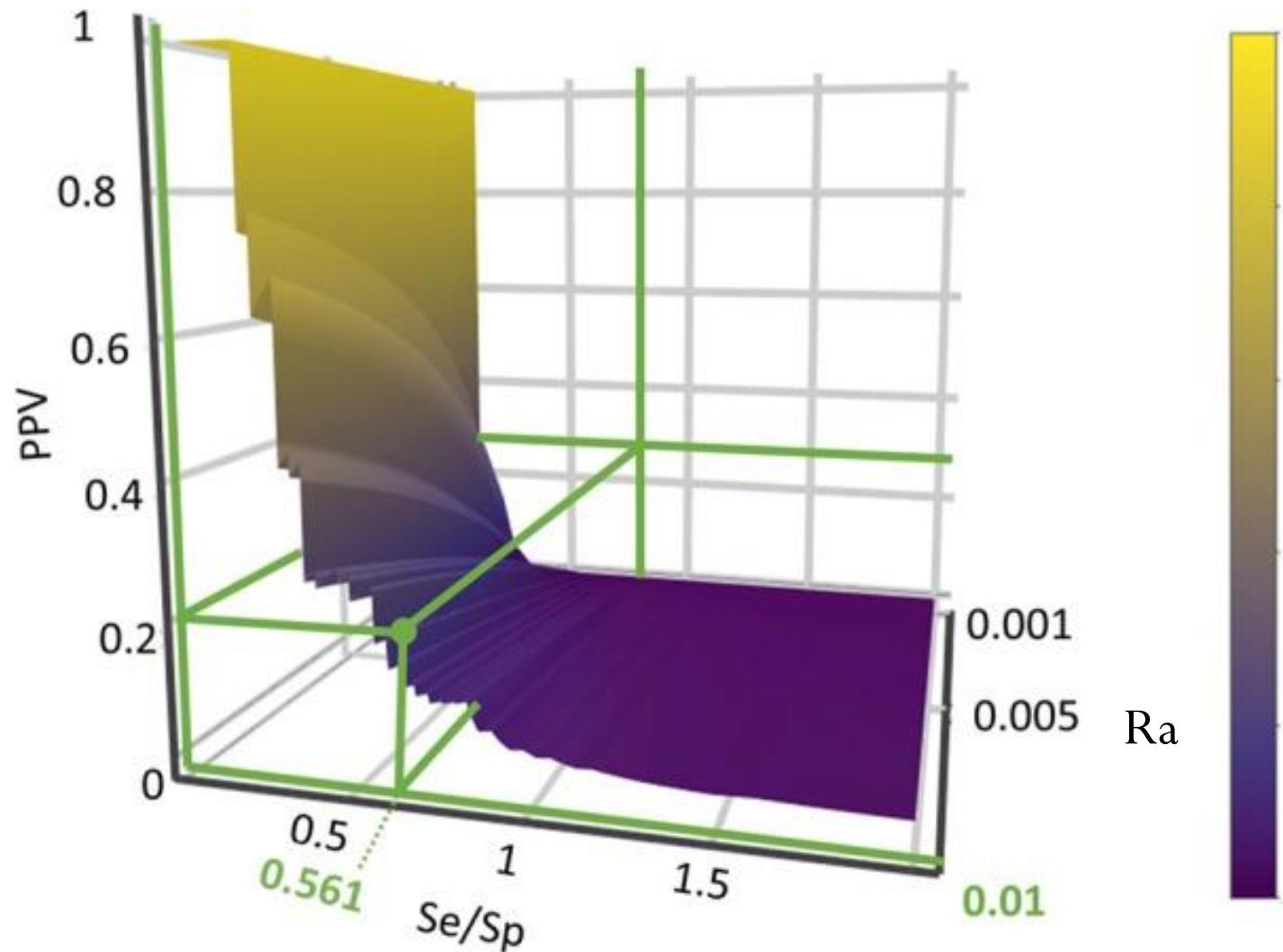
RIE; BEDROC: trabajar con  $\alpha R_a << 1$

Una consideración práctica:  
cuántos hits in silico deberé evaluar para obtener un  
hit confirmado (experimentalmente)?

$$PPV = \frac{Se \ Ra}{Se \ Ra + (1 - Sp)(1 - Ra)}$$

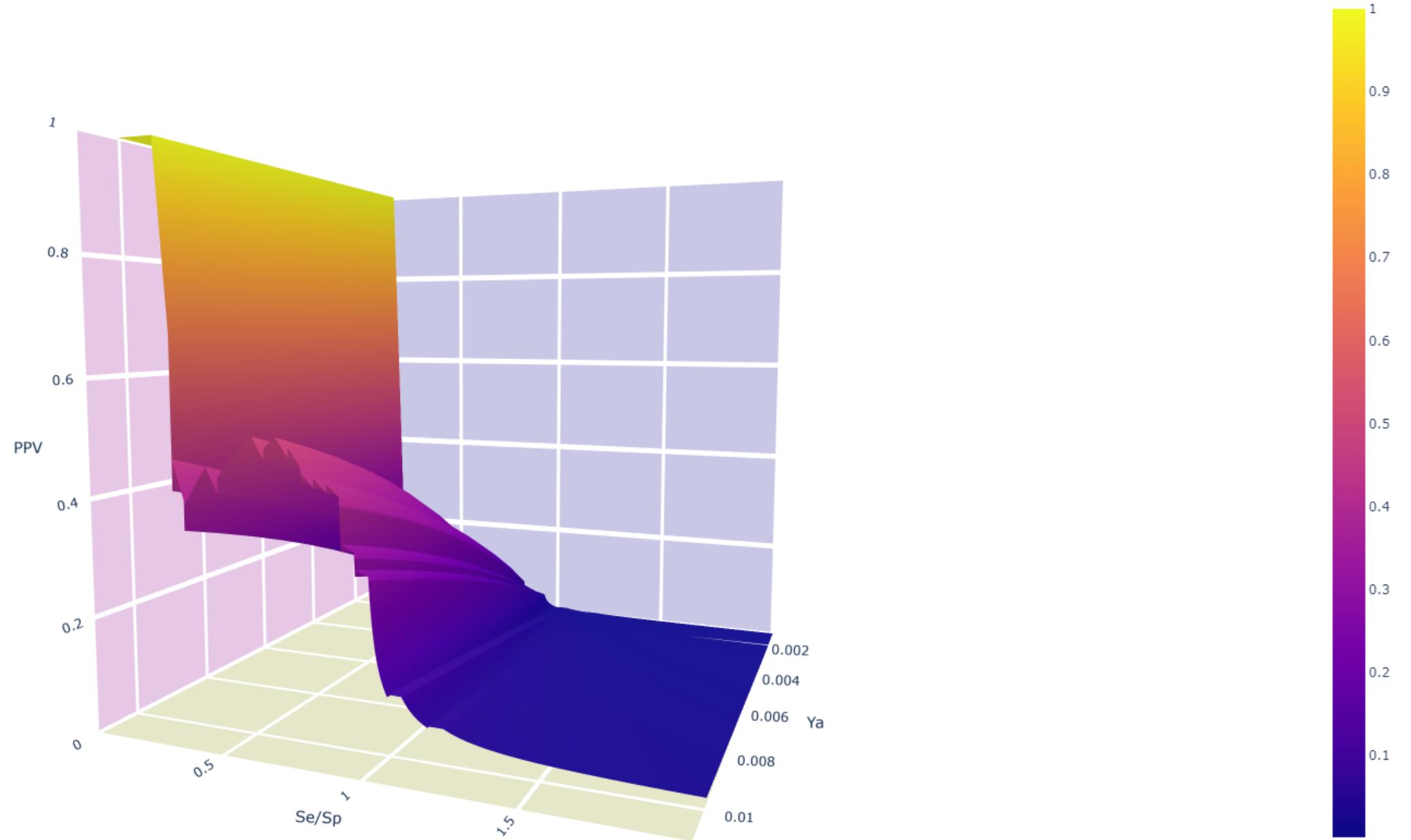
En un experimento de cribado real, Ra no se conoce a priori  
(suele variar, sin embargo, entre 0,01 y 0,005)

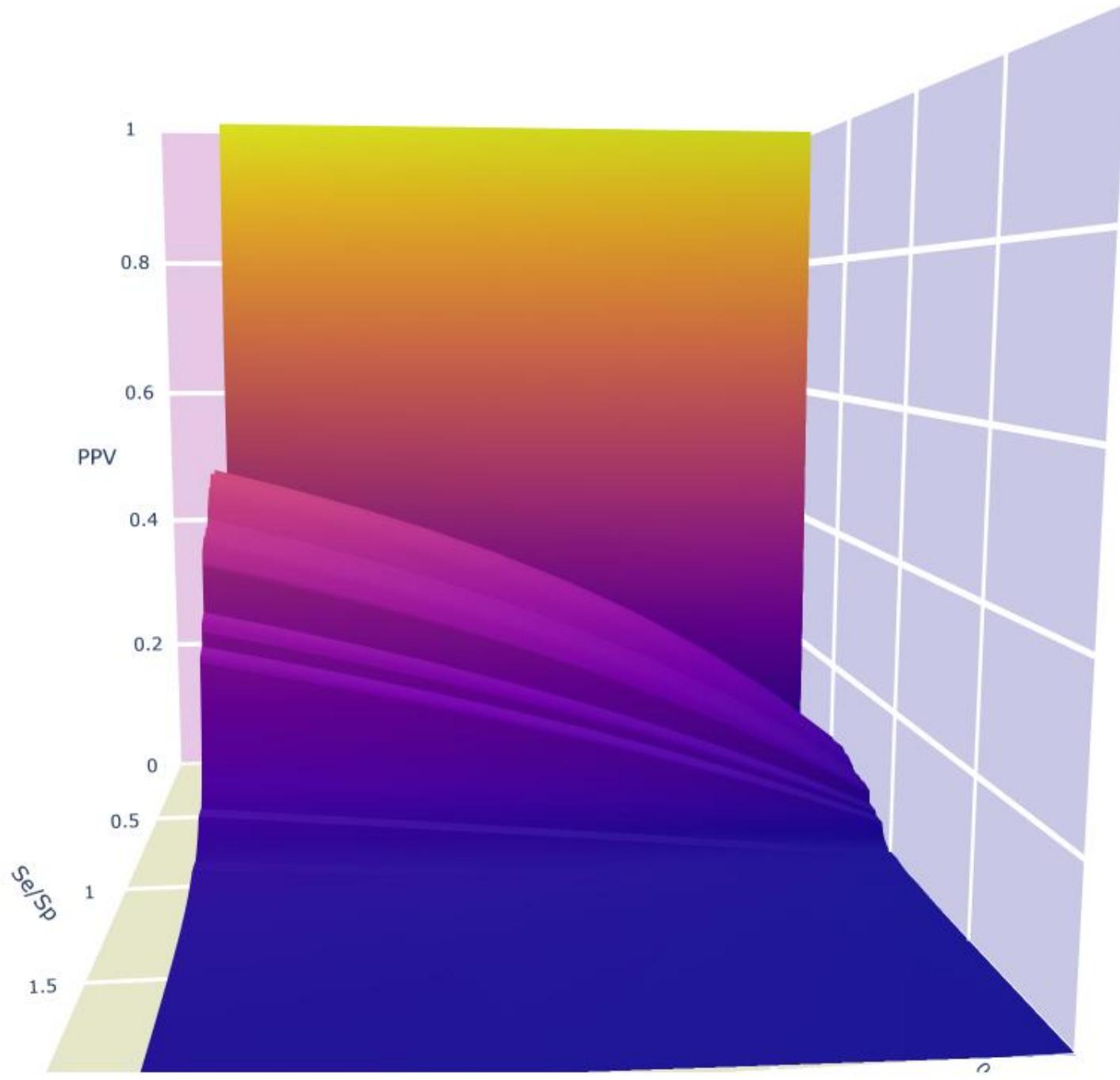
# Superficies PPV para optimizar el valor de corte del score

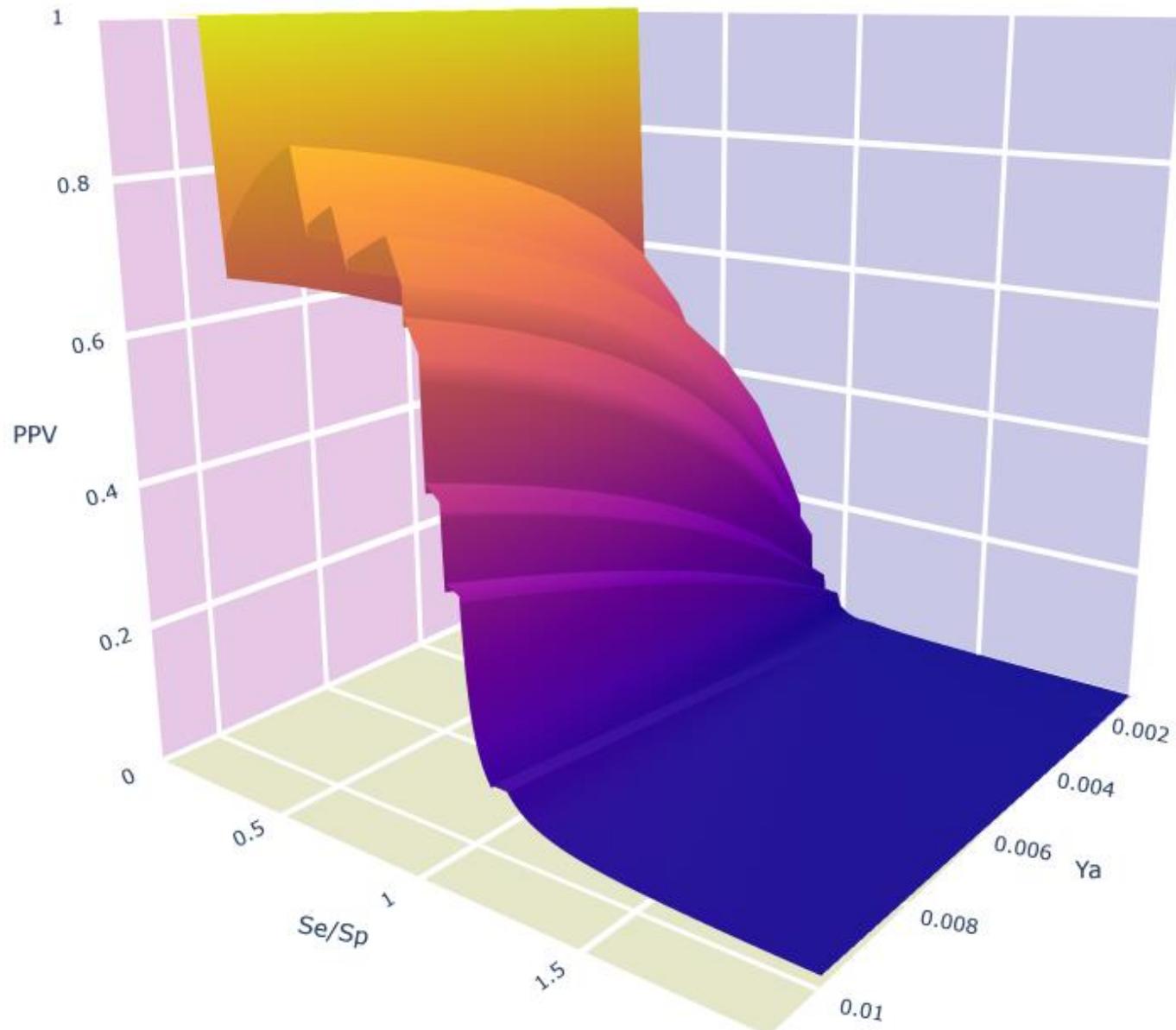


Alberca et al.  
Front. Cell. Infect. Microbiol.  
2018

Se asume que la relación Se/Sp observada en experimentos de cribado virtual retrospectivo adecuadamente diseñados se mantendrá aproximadamente en experimentos prospectivos

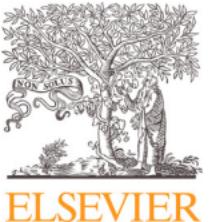








Artificial Intelligence in the Life Sciences 2 (2022) 100049



Contents lists available at [ScienceDirect](#)

## Artificial Intelligence in the Life Sciences

journal homepage: [www.elsevier.com/locate/ailsci](http://www.elsevier.com/locate/ailsci)



LIDeB Tools: A Latin American resource of freely available, open-source cheminformatics apps

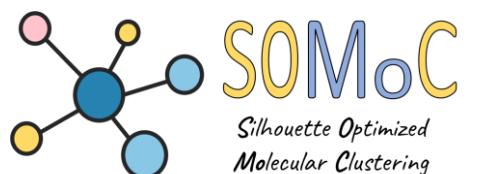
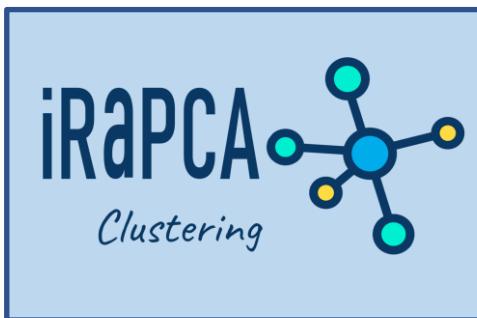
Denis N. Prada Gori, Lucas N. Alberca, Santiago Rodriguez, Juan I. Alice, Manuel A. Llanos, Carolina L. Bellera, Alan Talevi\*



*Laboratory of Bioactive Compounds Research and Development (LIDeB), Faculty of Exact Sciences, University of La Plata (UNLP), La Plata, Buenos Aires, Argentina*

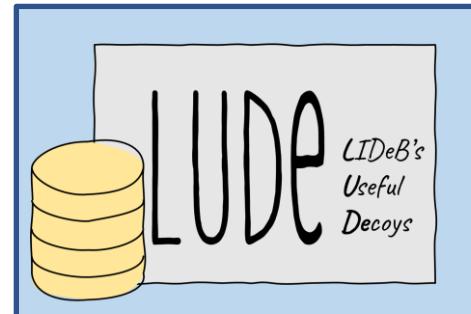


Clustering  
de pequeñas moléculas



**CHiCA**  
Comparative Hierarchical  
Clustering Algorithms

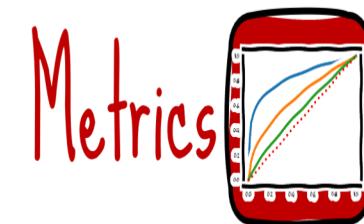
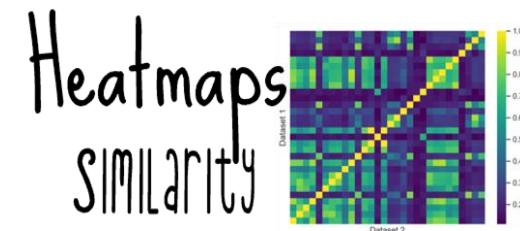
Generación de señuelos



Predictión de drogabilidad



Miscelánea

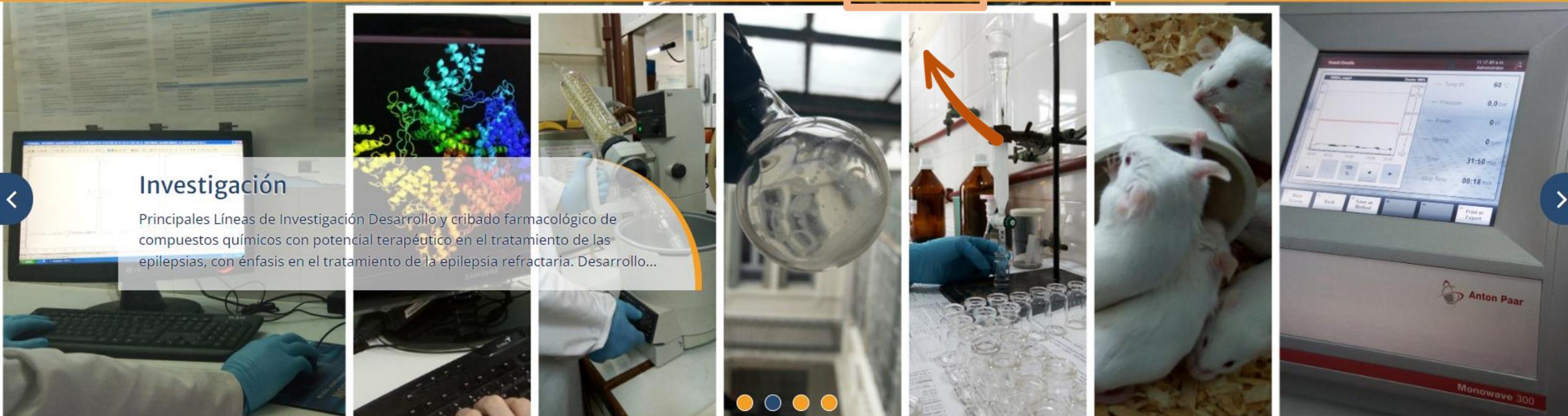


**LISTo**  
Standardization

Facultad de  
Ciencias Exactas  
Departamento de  
Ciencias BiológicasUNIVERSIDAD  
NACIONAL  
DE LA PLATA

# Laboratorio de Investigación y Desarrollo de Bioactivos

Facultad de Ciencias Exactas | Universidad Nacional de La Plata

[Inicio](#)[Institucional](#) ▾[Investigación](#) ▾[Docencia](#)[Transferencia](#)[LIDeB tools](#)[Novedades](#)[Contacto](#)

## Investigación

Principales Líneas de Investigación Desarrollo y cribado farmacológico de compuestos químicos con potencial terapéutico en el tratamiento de las epilepsias, con énfasis en el tratamiento de la epilepsia refractaria. Desarrollo...

## ¿Quienes Somos?

Laboratorio de Investigación  
y Desarrollo de BioactivosFacultad de Ciencias Exactas  
Universidad Nacional de La Plata

## Investigación



Principales Líneas de Investigación Desarrollo v

## Publicaciones



Publicaciones últimos 5 años 2021 Llanos M..



## >About Us

We are a drug discovery team with an interest in the development of publicly available open-source customizable cheminformatics tools to be used in computer-assisted drug discovery. We belong to the Laboratory of Bioactive Research and Development (LiDeB) of the National University of La Plata (UNLP), Argentina. Our research group is focused on computer-guided drug repurposing and rational discovery of new drug candidates to treat epilepsy and neglected tropical diseases.

⚡ Coming soon more webapps! 🎉

## Heatmaps- Similarity



### About this App

Build a Heatmap of molecular similarity. These plots of inter-molecular similarity (computed as Tanimoto similarity coefficient using Morgan fingerprints and other molecular fingerprinting systems) allow for a fast, visual inspection of the molecular diversity of the datasets, and also preliminary detection of clusters within a dataset. The resulting plots are downloadable as .png files through a simple right click on your mouse!

[App](#) [Publication](#)

## iRaPCA-Clustering



### About this App

iRaPCA Clustering is based on an iterative combination of the random subspace approach (feature bagging), dimensionality reduction through Principal Component Analysis (PCA) and the k-means algorithm. The optimal number of clusters k and the best subset of descriptors are selected from plots of silhouette coefficient against different k values and subsets. Different validation metrics can be downloaded once the process is finished. A number of graphs may be built and readily downloaded through a simple click.

[App](#) [Publication](#)

## LuDe



## SOMoC Clustering



- » Importante financiamiento de Apoyo Dravet para el LiDeB
- » Nuevos avances en la lucha contra la toxoplasmosis
- » LiDeB tools
- » 3rd International online congress on Dravet Syndrome and Refractory Epilepsy
- » Seminario en IPMontevideo

## Líneas de Investigación



## Sitios de Interés

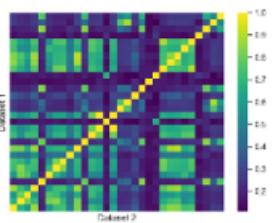
- » Facultad de Ciencias Exactas
- » Universidad Nacional de La Plata
- » CONICET
- » INDRE Network
- » Departamento Ciencias Biológicas

## Seguinos en:



## Heatmaps- Similarity

# Heatmaps SIMILARITY



### About this App

Build a Heatmap of molecular similarity. These plots of inter-molecular similarity (computed as Tanimoto similarity coefficient using Morgan fingerprints and other molecular fingerprinting systems) allow for a fast, visual inspection of the molecular diversity of the datasets, and also preliminary detection of clusters within a dataset. The resulting plots are downloadable as .png files through a simple right click on your mouse!

[App](#)[Publication](#)

web app

## iRaPCA-Clustering



### About this App

iRaPCA Clustering is based on an iterative combination of the random subspace approach (feature bagging), dimensionality reduction through Principal Component Analysis (PCA) and the k-means algorithm. The optimal number of clusters k and the best subset of descriptors are selected from plots of silhouette coefficient against different k values and subsets. Different validation metrics can be downloaded once the process is finished. A number of graphs may be built and readily downloaded through a simple click.

[App](#)[Publication](#)

Artículo publicado

Seguinos en:



**Upload your SMILES**

Upload a TXT file with one SMILES per line

Drag and drop file here  
Limit 200MB per file • TXT

[Browse files](#)

[Example TXT input file](#)

Check to change the default configuration

**Physicochemical features limits**

- Molecular weight: 5
- LogP: 0.50
- Rotable bonds: 1
- Num of H Acceptors: 1
- Num of H Donors: 1

### eB Tools - LUDe

LUDe (LIDeB's Useful Decoys) is a WebApp that generates, from a set of active compounds, decoys (putative inactive compounds) which can be used to retrospectively validate virtual screening tools/protocols. Decoys are molecules that have not been tested against a molecular target of interest but due to their structural features are presumably not prone to bind the target with the known active compounds in relation to certain general physicochemical properties (e.g., molecular weight, log P, and others) but are topologically different from the query compounds. LUDe is conceptually similar to the Directory of Useful Decoys enhanced, but additional filters have been serially implemented to assure the topological dissimilarity between the decoys and the query App, decoys are obtained through four sequential steps:

- Choosing molecules with similar physicochemical properties of the input active molecules in a curated CHEMBL database.
- Ranking the selected molecules by dissimilarity against each individual input molecule.
- Simply selecting a desired number of decoys for each individual input molecule.
- Ranking the selected molecules by the dissimilarity against all the input molecules (set of active compounds used as query). The decoys will have low Tanimoto similarity with the input compounds, and also low Maximum Common Substructure (MCS) ratio and distinctive molecular frameworks (Murcko scaffold). All in all, these three serial filters assure that the decoys will be chemically similar to the active compounds but dissimilar to the remaining active compounds is also checked. Finally, you can download a file with your decoys.

The workflow summarizes the steps performed by this method:

```

graph LR
    Input[Input: .txt file with one SMILES per line] --> Standardization[Standardization? ... Yes, we do!]
    Standardization --> Search[Searching for compounds with the selected properties in CHEMBL27 Database]
    Search --> Properties[MOLECULAR WEIGHT LogP ± 0.5 N° OF ROTATABLE BONDS N° OF H DONORS N° OF H ACCEPTORS ± 2]
    Properties --> Decoys[LUDe LIDeB's Useful Decoys]
    Decoys --> Similarity[SUBSET OF COMPOUNDS WITH SIMILAR PROPERTIES TO EACH ACTIVE]
    Similarity --> Settings[SMILES TANIMOTO COEFFICIENT 0.2]
    Settings --> MaxDecoys[Maximum number of decoys: Maximum number of decoys by active compound 100]
    MaxDecoys --> Outputs[Outputs: .txt file with decoys .txt file with settings]
    
```

**Contact Us**

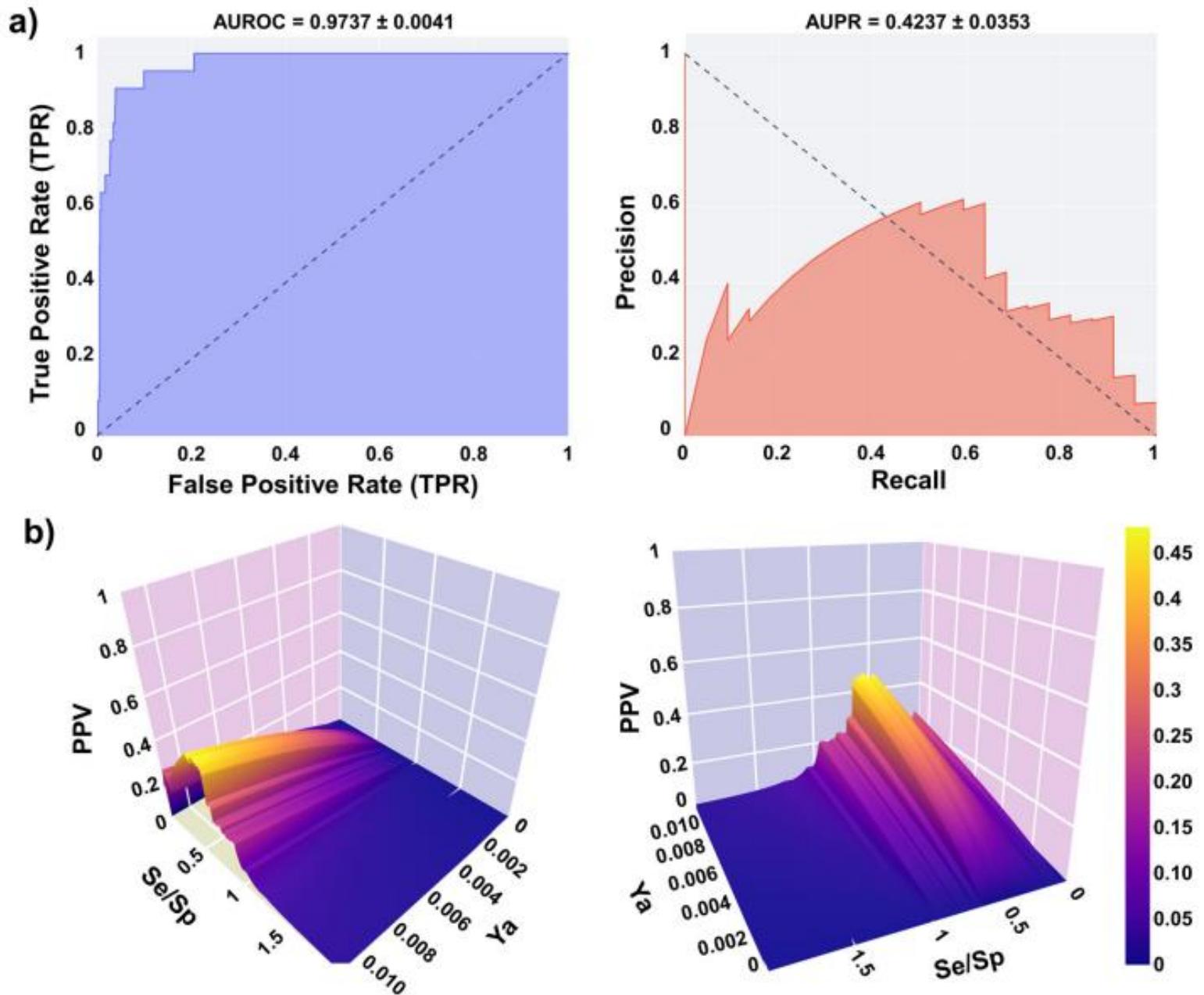
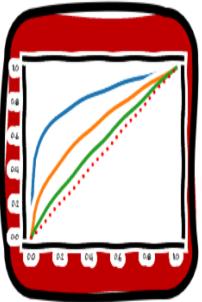
If you are looking to contact us, please [Email](#) or [Twitter](#)

**Streamlit**

el código también está disponible en [github.com](https://github.com/lideb)

<https://github.com/lideb>

# Metrics



# Generación de señuelos

estos no son buenos decoys



estos tampoco son buenos decoys



los decoys son compuestos que en cierto sentido muy general son parecidos a los compuestos activos que se conocen contra determinado blanco farmacológico, pero que son topológicamente distintos



## Directory of Useful Decoys, Enhanced (DUD-E): Better Ligands and Decoys for Better Benchmarking

Michael M. Mysinger,<sup>†</sup> Michael Carchia,<sup>†</sup> John. J. Irwin,<sup>\*,†</sup> and Brian K. Shoichet<sup>\*,†</sup>

<sup>†</sup>Department of Pharmaceutical Chemistry, University of California San Francisco, San Francisco, California 94158-2330, United States

## DeepCoy

Structural bioinformatics

### Generating property-matched decoy molecules using deep learning

Fergus Imrie <sup>1</sup>, Anthony R. Bradley<sup>2</sup> and Charlotte M. Deane <sup>1,\*</sup>

## DUD-E

Article  
[pubs.acs.org/jmc](https://pubs.acs.org/jmc)

*Bioinformatics*, 37(15), 2021, 2134–2141

doi: 10.1093/bioinformatics/btab080

Advance Access Publication Date: 3 February 2021

Original Paper



JOURNAL OF  
CHEMICAL INFORMATION  
AND MODELING

[pubs.acs.org/jcim](https://pubs.acs.org/jcim)

## DUD-Z

Article

### Property-Unmatched Decoys in Docking Benchmarks

Reed M. Stein, Ying Yang, Trent E. Baliaus, Matt J. O'Meara, Jiankun Lyu, Jennifer Young, Khanh Tang, Brian K. Shoichet,<sup>\*</sup> and John J. Irwin<sup>\*</sup>

Cite This: *J. Chem. Inf. Model.* 2021, 61, 699–714

Read Online

## DecoyFinder

BIOINFORMATICS APPLICATIONS NOTE

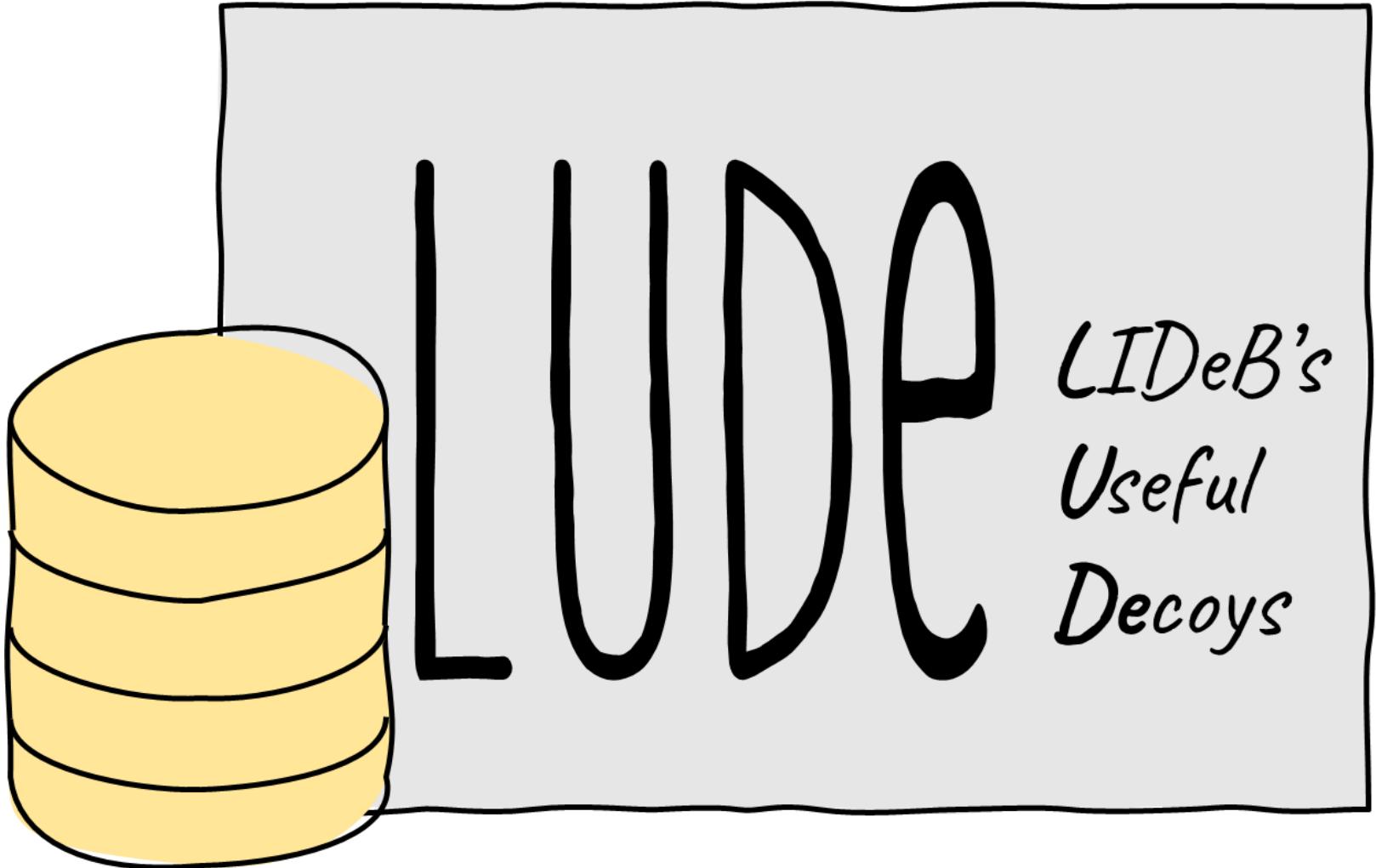
Vol. 28 no. 12 2012, pages 1661–1662  
doi:10.1093/bioinformatics/bts249

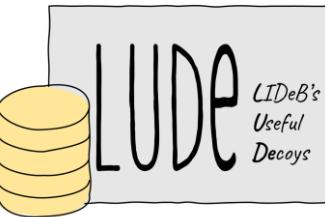
Structural bioinformatics

Advance Access publication April 26, 2012

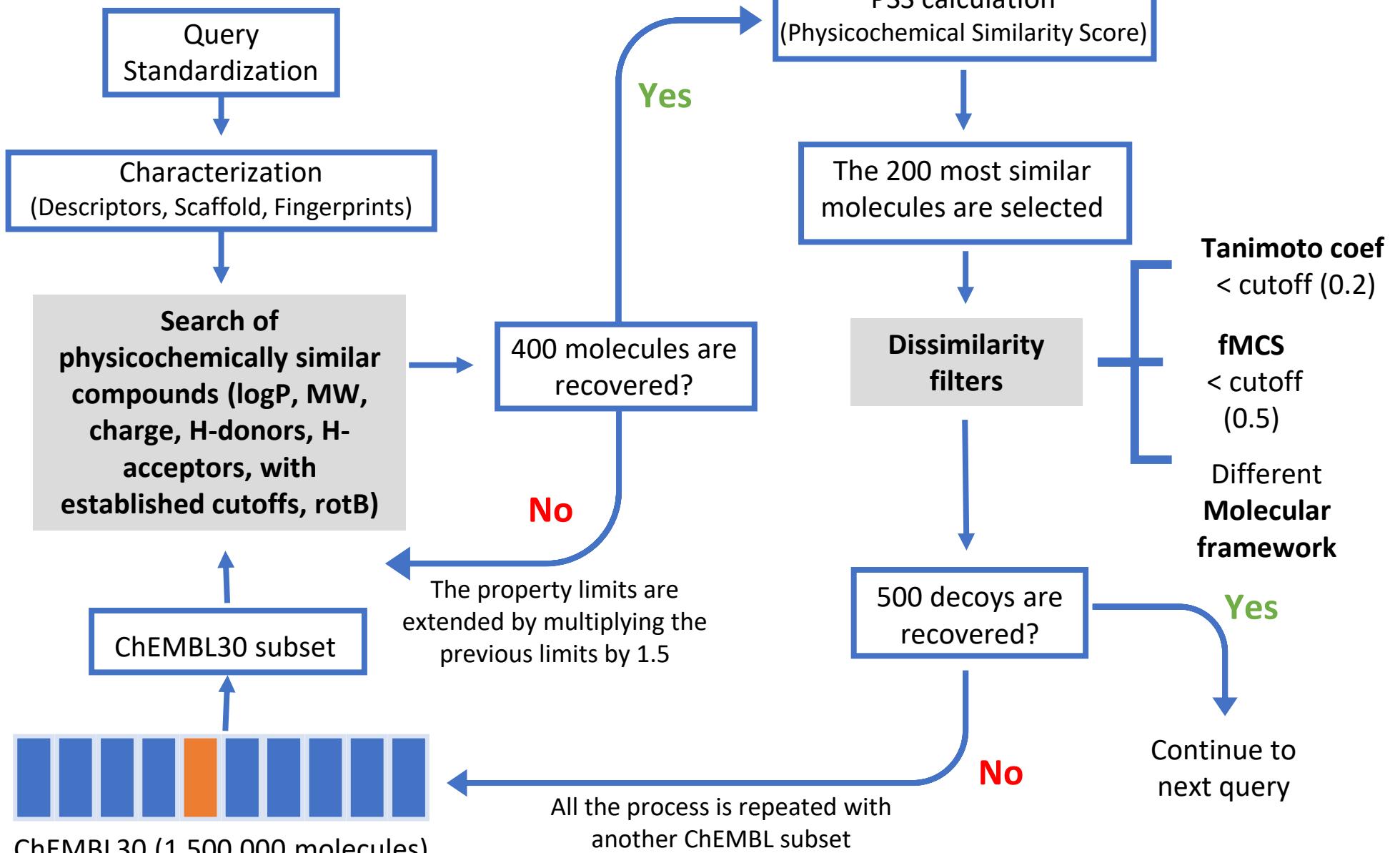
### DecoyFinder: an easy-to-use python GUI application for building target-specific decoy sets

Adrià Cereto-Massagué<sup>1</sup>, Laura Guasch<sup>1</sup>, Cristina Valls<sup>1</sup>, Miquel Mulero<sup>1</sup>, Gerard Pujadas<sup>1,2</sup> and Santiago Garcia-Vallvé<sup>1,2,\*</sup>





## Workflow for each query



## Workflow with all decoys and queries

Similarity between decoys from each query and all other queries is computed

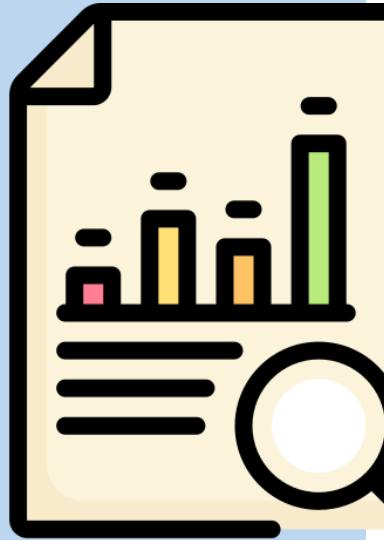
Decoys with  $T_c >$  selected cutoff are removed

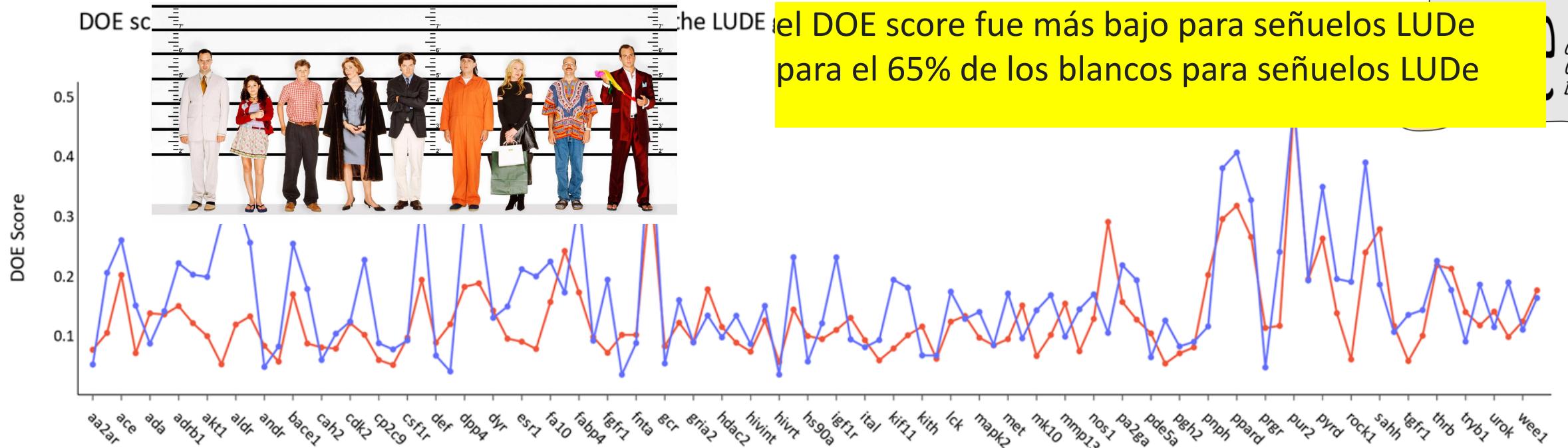
50 decoys are randomly selected

# BENCHMARKING

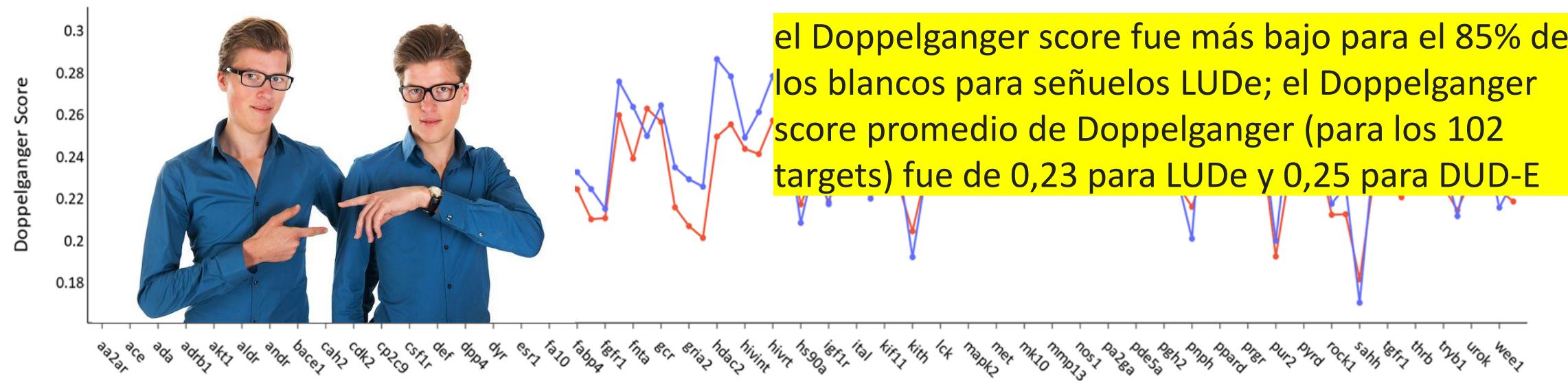
## Validación de decoys

- **Sets de datos:**
  - 102 blancos farmacológicos extraídos de DUD-E
    - 50 decoys por compuesto activo
- **Métricas:**
  - DOE score: para cada query (activo) se calcula la distancia en el espacio de las 6 matching properties normalizadas, contra todos los demás (activos y señuelos). Se etiquetan como 1 a los activos y como 0 a los inactivos y se construye la curva ROC. Luego se calcula el promedio de  $ROC - 0,5$  para todos los activos.
  - Doppelganger score. Máxima similitud de Tanimoto un par cualquiera de los queries y los decoys.

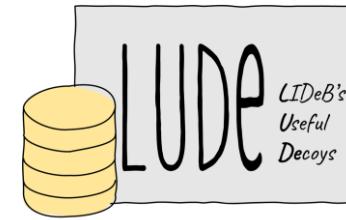




Doppelganger scores of the original DUD-E set (blue) compared to the LUDe generated decoys (red)



# Ejemplo: Dataset fabp4 (Fatty Acid Binding Protein 4)



## OUTPUT 2 Generated\_decoys.csv

SMILE	Query
Clc1ccc(C(CCCOc2ccc(Br)cc2)Cn2ccnc2)c(Cl)c1	Query_1
CCOC(=O)CC[C@H](O[C@H]1O[C@H]2O[C@H]3(C)CC[C@H]4[C@H](C)CC[C@H](C[C@H]1C)[C@]42OO3)c1cccc1	Query_1
CC(C)(F)c1noc(N2CCC(COC3CC=C(c4ccc(S(C)(=O)=O)cc4)CC3)CC2)n1	Query_1
COc1cc(c2ccc(COc3cc([C@@H](CC(=O)[O-])C4CC4)ccn3)cc2[C@@H]2CCCC2(C)C)c(F)cn1	Query_1
O=C([O-])CCc1cccc1/C=C1/C2CCC(O2)C1c1nc(C(=O)NCCCCC2CCCCC2)co1	Query_1
O=C1[C@H](N2C(=O)c3cccc3C2=O)CN1[C@H](COCc1cccc1)c1ccc(OCc2cccc2)cc1	Query_1
CCCc1cnc(N2CCC(C3Cc4cc(C5CCN(S(=O)(=O)CCC)CC5)ccc4O3)CC2)nc1	Query_1
Cc1ccc(C(C)C)c(SC2=C([O-])O[C@@](CCc3cccc3)(c3cccc3)CC2=O)c1	Query_1
O=C(CC(CC(=O)c1ccc(Br)cc1)c1ccnc1)c1ccc(Br)cc1	Query_1
CN(Cc1ccnc(N2CCCC2)c1)C(=O)c1cnc(Oc2ccc3c(c2)CCC(c2cccc2)O3)s1	Query_1
Cc1ccc(c2nc3sc(C)nn3c2COCc2cn(Cc3ccc(F)cc3)nn2)cc1	Query_1
CC(C)c1ccc(C2CC(=O)c3ccc(OCc4cn(Cc5ccc(Br)cc5)nn4)cc3O2)c1	Query_1
-----	Query_X

## OUTPUT 3 Decoys\_settings.csv

Decoys generated at:	8/8/2022
Physicochemical features limits:	
MW	+/- 20
logP	+/- 0.5
Num_rotatable_bonds	+/- 1
Num_H_acceptors	+/- 1
Num_H_donors	+/- 1
Topological features ---> Dissimilarity conditions	
Morgan Fingerprints	
Fingerprint radio:	2
Fingerprint lenght:	1024
Limit of fraction of the Maximum Common Substructure:	0.5
Decoys with different framework:	True
Max Tc similarity between decoys and other actives:	0.2
Max number of decoys by loaded molecule:	50





**NO USAR** señuelos generados utilizando criterios de exclusión basados en similitud molecular para evaluar métodos basados en similitud molecular.

Los creadores de DUD-E indican que sus señuelos deben usarse exclusivamente para evaluar el desempeño de protocolos de docking molecular.

# Clustering

## Clustering: aplicaciones

Muestreo representativo de quimiotecas / bases de datos químicas.  
Por ejemplo, para particionar representativamente un dataset en training,  
test, validation sets

Muestreo representativo de hits in silico

Aprendizaje automático no supervisado

## Clustering: clasificación

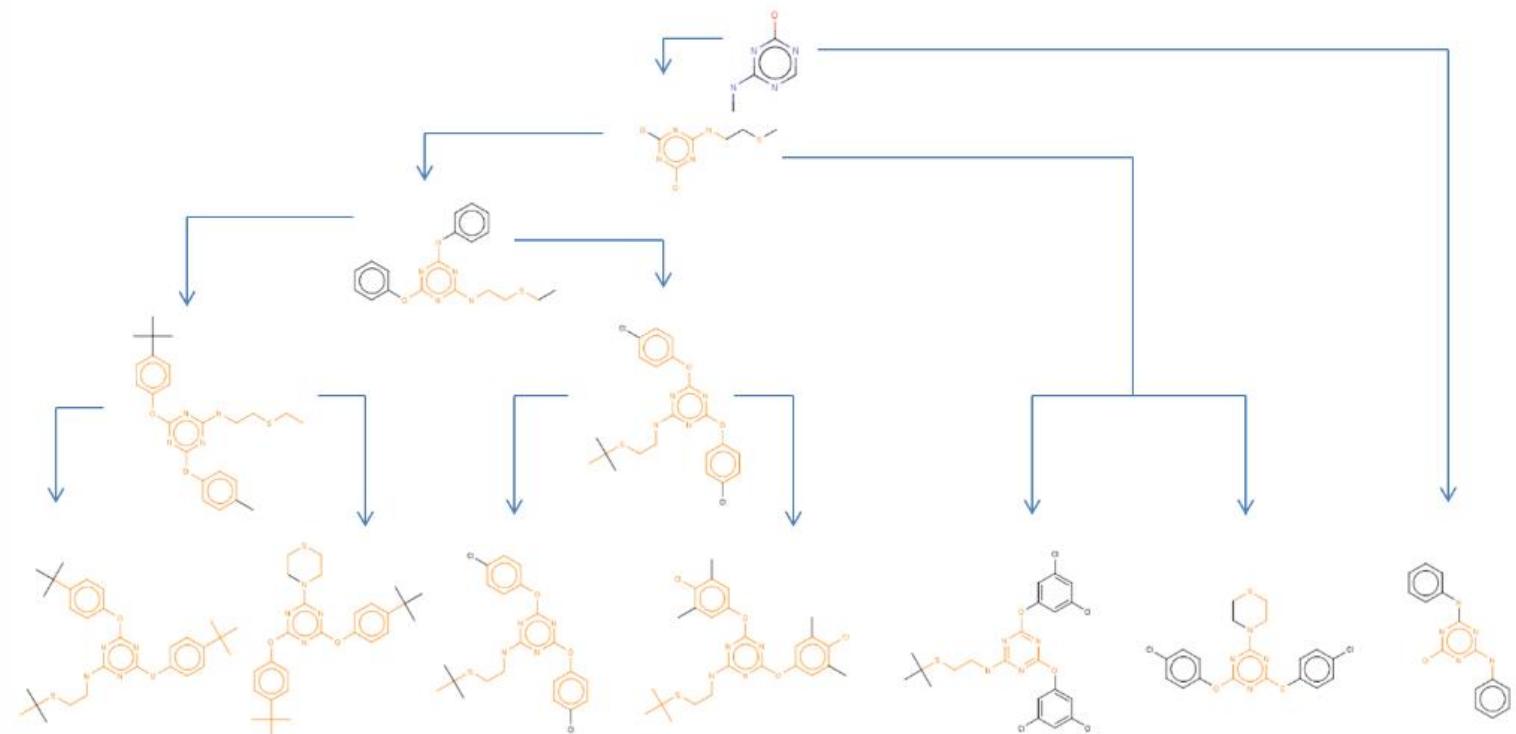
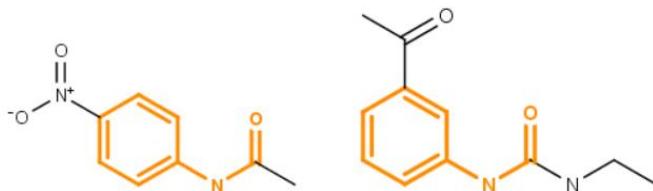
Jerárquico (métodos aglomerativos o divisivos): Single linkage, complete linkage, average linkage, Chemaxon's Library MCS.

No jerárquico: Esferas de exclusión, K-means.

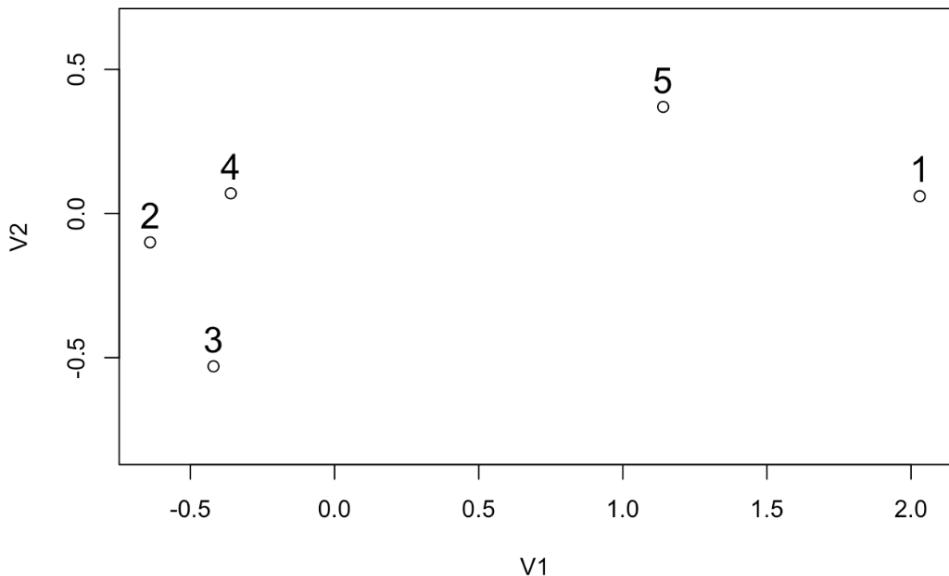
Híbridos: Library MCS + k-means.

## Library MCS

Se basa en agrupar las moléculas en función de su máxima subestructura común. El usuario puede definir el tamaño mínimo de la MCS y si los anillos se considerarán como tales o puede romperse. Se inicia agrupando las estructuras que presentan una MCS de mayor tamaño.



## Single y complete linkage

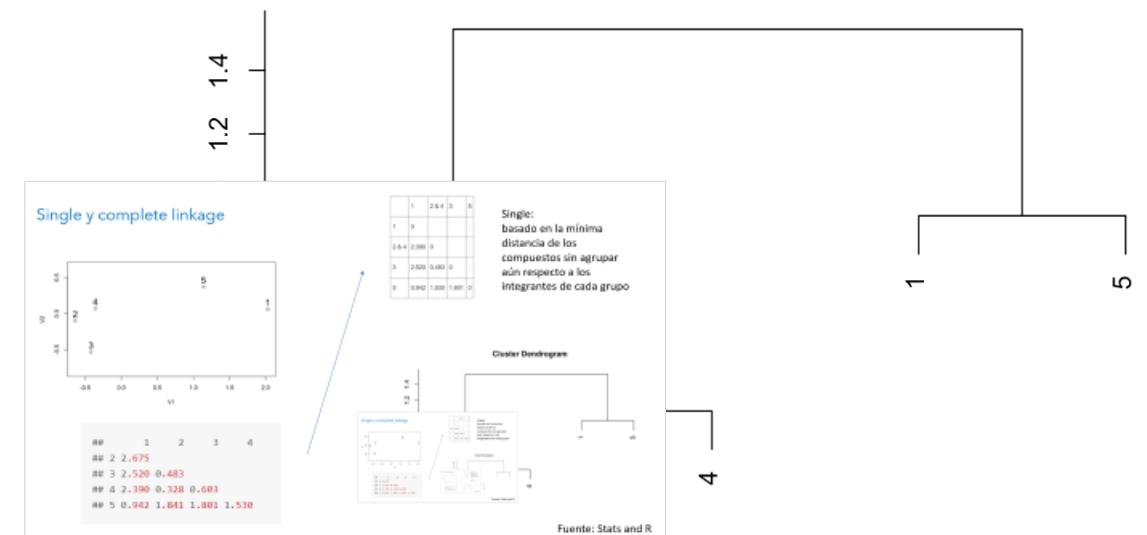


```
##      1   2   3   4  
## 2 2.675  
## 3 2.520 0.483  
## 4 2.390 0.328 0.603  
## 5 0.942 1.841 1.801 1.530
```

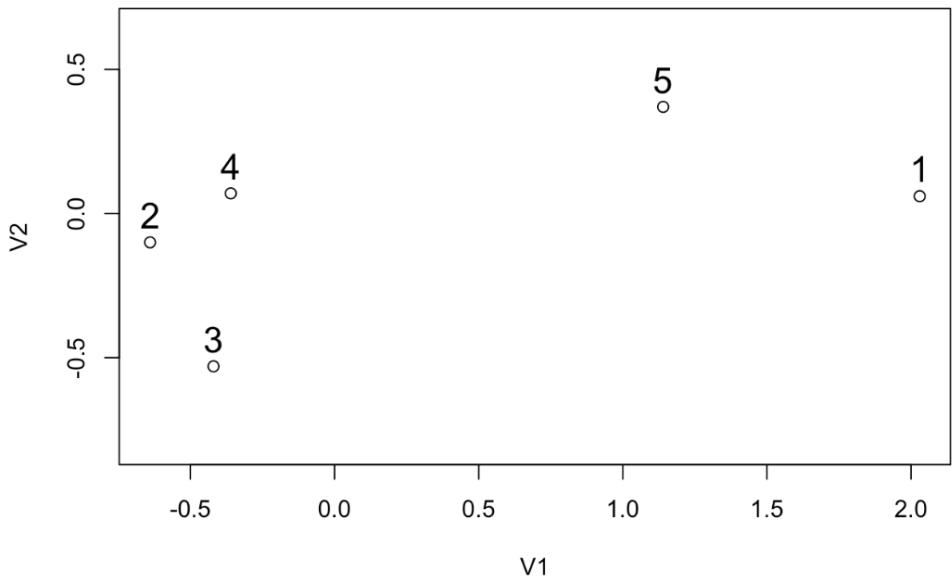
	1	2 & 4	3	5
1	0			
2 & 4	2.390	0		
3	2.520	0.483	0	
5	0.942	1.530	1.801	0

**Single:**  
basado en la mínima distancia de los compuestos sin agrupar aún respecto a los integrantes de cada grupo

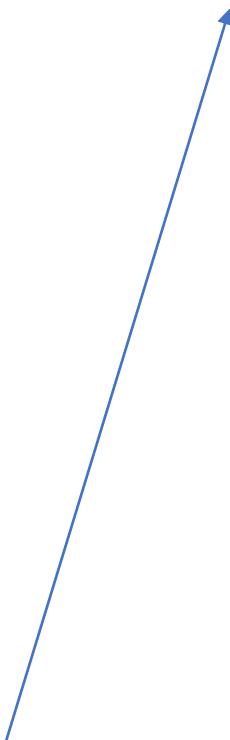
Cluster Dendrogram



## Single y complete linkage



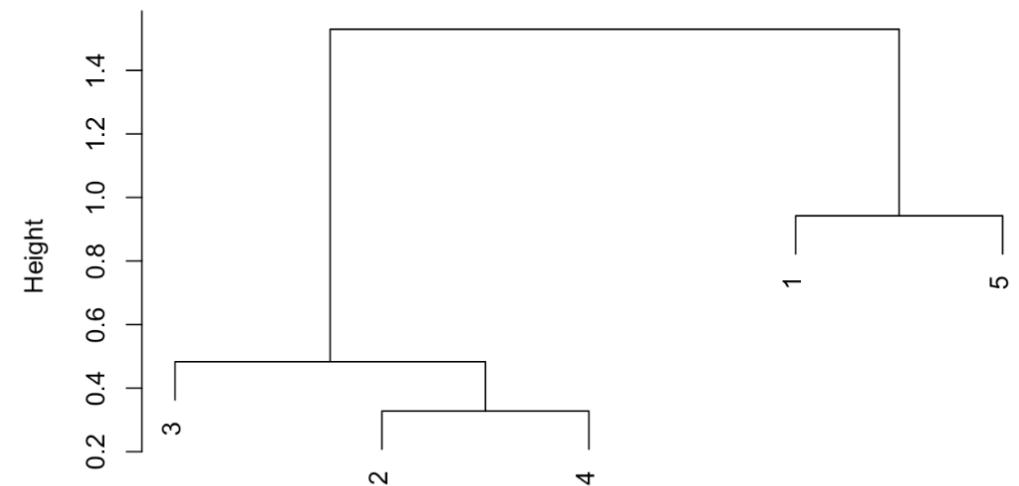
```
##      1   2   3   4
## 2 2.675
## 3 2.520 0.483
## 4 2.390 0.328 0.603
## 5 0.942 1.841 1.801 1.530
```



	1	2 & 4	3	5
1	0			
2 & 4	2.675	0		
3	2.520	0.603	0	
5	0.942	1.841	1.801	0

Complete:  
basado en la máxima  
distancia de los  
compuestos sin agrupar  
aún respecto a los  
integrantes de cada grupo

Cluster Dendrogram



# Butina

## 1. Calculate fingerprints

- ▲ 1011011...
- 1100010...
- 1010001...
- 1110110...
- 1011001...
- ▲ 1011111...
- 0100010...
- ▲ 1001111...
- 1010000...

## 2. Tanimoto similarity matrix

threshold = .7

	▲	●	■	●	■	▲	●	▲	■
▲	.3	.7	.4	.7	.9	.2	.5	.7	
●		.2	.8	.4	.6	.9	.7	.7	
■			.4	.9	.7	.2	.6	.9	
●				.5	.6	.9	.7	.6	
■					.7	.3	.7	.8	
▲						.4	.9	.6	
●							.5	.2	
▲								.7	

⇒ sort molecules by #neighbors  
(molecules with similarity ≥ .7)

	#neighbors	#neighbors
▲	4	5
●	4	5
■	4	5
●	3	4
■	5	4
▲	4	4
●	2	4
▲	5	3
■	5	2

sort  
⇒

## 3. Clustering

Flag cluster members, start from top of the list:

1. Iteration: molecule ■

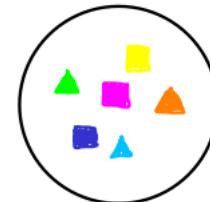
Neighbors:

	▲	●	■	●	■	▲	●	▲	■
▲	.3	.7	.4	.7	.9	.2	.5	.7	
●		.2	.8	.4	.6	.9	.7	.7	
■			.4	.9	.7	.2	.6	.9	
●				.5	.6	.9	.7	.6	
■					.7	.3	.7	.8	
▲						.4	.9	.6	
●							.5	.2	
▲								.7	

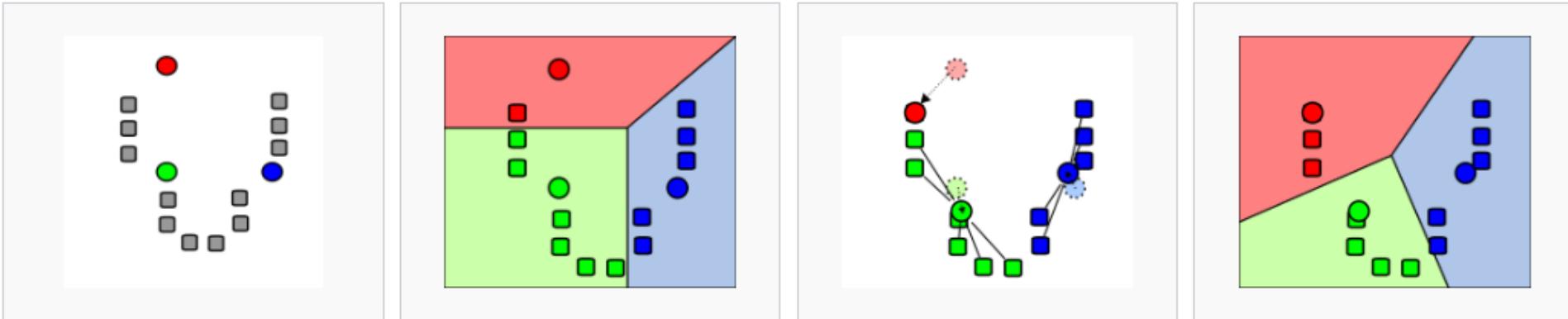
Flag neighbors (and centroid):

	▲	●	■	●	■	▲	●	▲	■
▲	.3	.7	.4	.7	.9	.2	.5	.7	
●		.2	.8	.4	.6	.9	.7	.7	
■			.4	.9	.7	.2	.6	.9	
●				.5	.6	.9	.7	.6	
■					.7	.3	.7	.8	
▲						.4	.9	.6	
●							.5	.2	
▲								.7	

Build new cluster:



# K-means



1)  $k$  centroides iniciales (en este caso  $k=3$ ) son generados aleatoriamente dentro de un conjunto de datos (mostrados en color).

2)  $k$  grupos son generados asociándole el punto con la media más cercana. La partición aquí representa el [diagrama de Voronoi](#) generado por los centroides.

3) El [centroide](#) de cada uno de los  $k$  grupos se recalcula.

4) Pasos 2 y 3 se repiten hasta que se logre la convergencia.

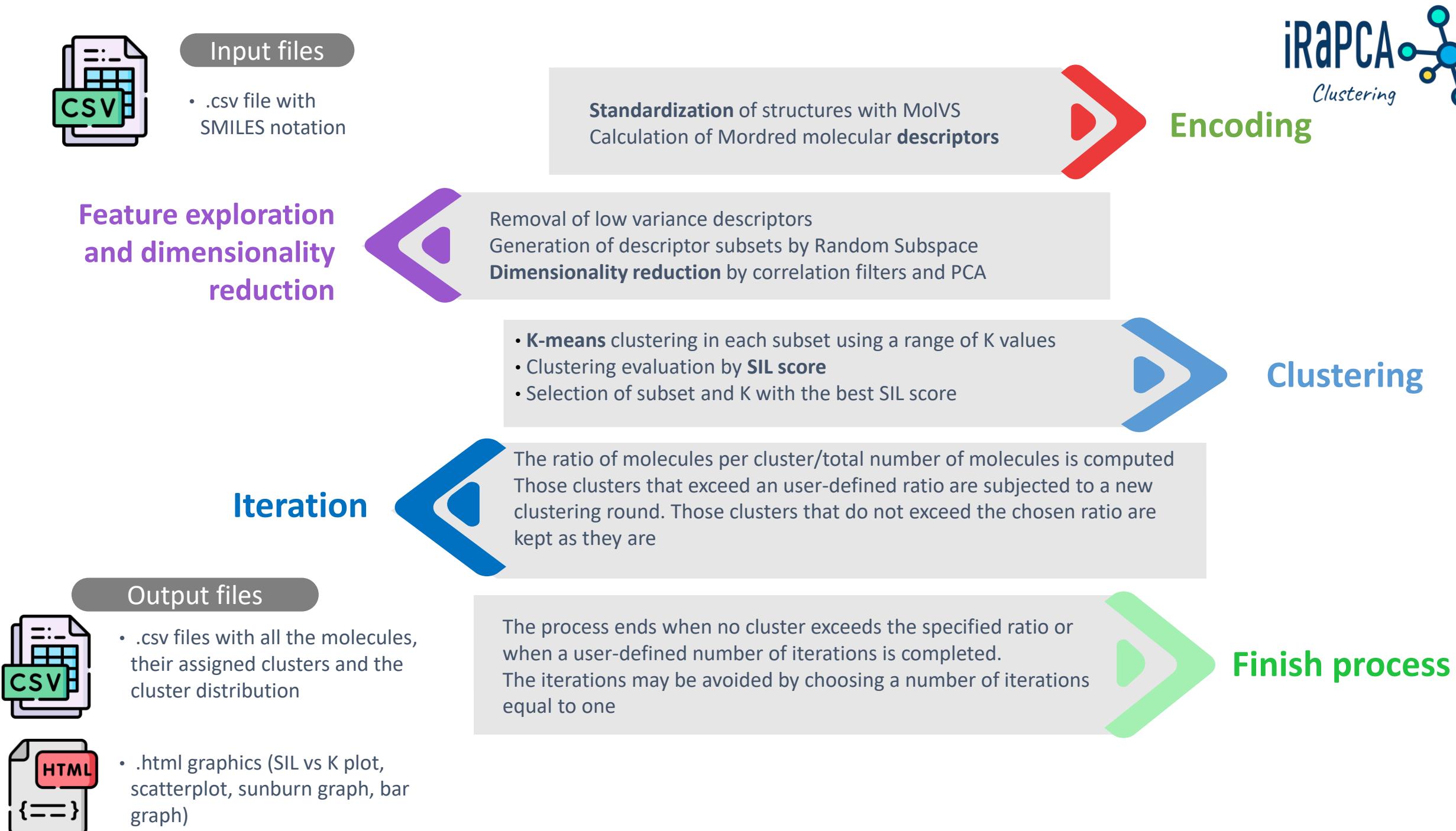
# iRAPCA

*Clustering*



## Iterative Random Subspace PCA clustering (iRaPCA)

- iRaPCA clustering se basa en una combinación iterativa de subespacios aleatorios (feature bagging), reducción de dimensionalidad mediante Análisis de Componentes principales (PCA) y el algoritmo K-means algorithm
- El valor medio de silhouette se utiliza habitualmente como métrica de validación del clustering. Para cada elemento del conjunto de datos,  $s(i)$  puede ser calculado como  $b(i)-a(i)/\max\{a(i),b(i)\}$ , donde  $a(i)$  es la distancia media entre  $i$  y todos los demás puntos de datos en el mismo cluster y  $b(i)$  es la distancia media más pequeña de  $i$  a todos los puntos en cualquier otro grupo.

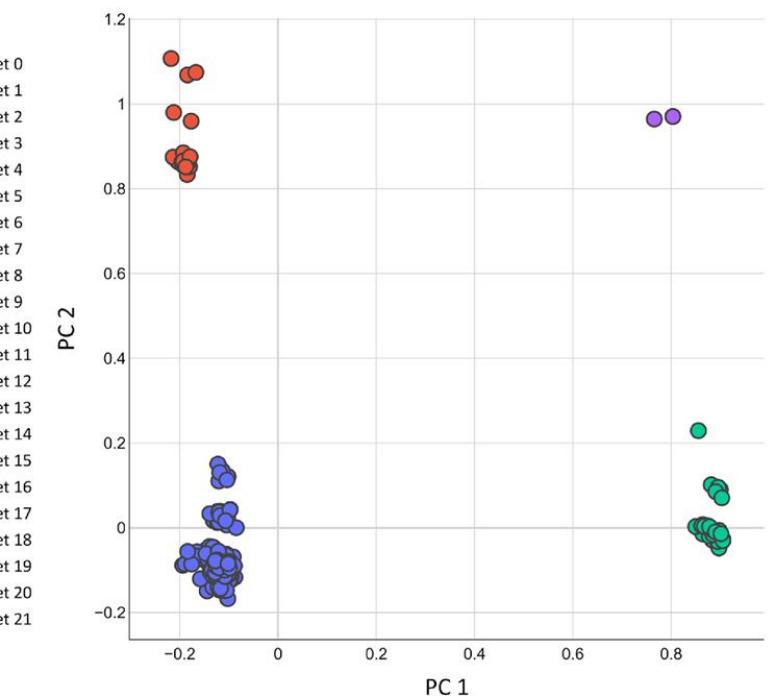
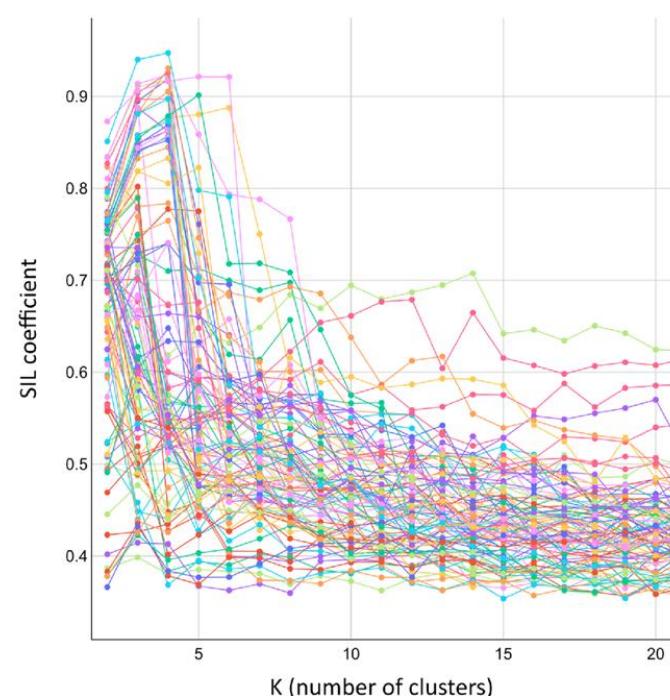


# Ejemplo: Dataset focal\_adhesion

## INPUT

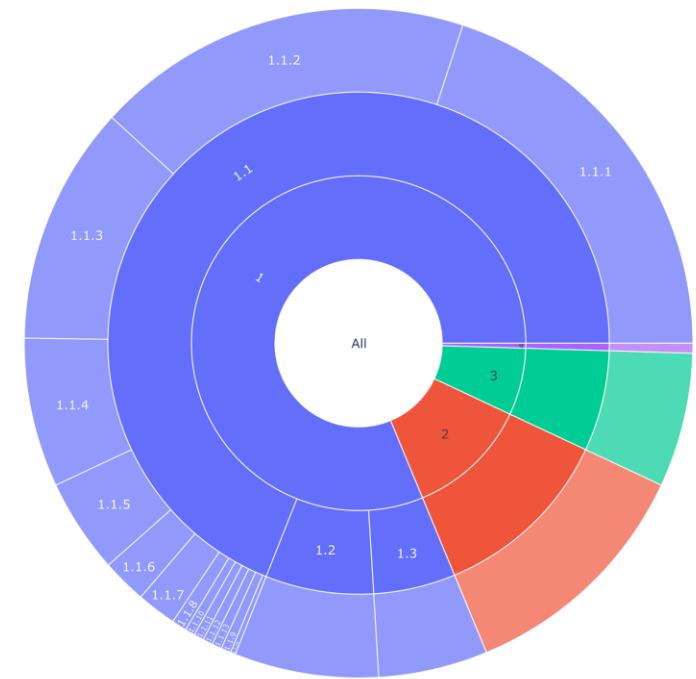


Archivo .csv con **416 moléculas**  
en anotación SMILES



## OUTPUT 1

Plots interactivos:



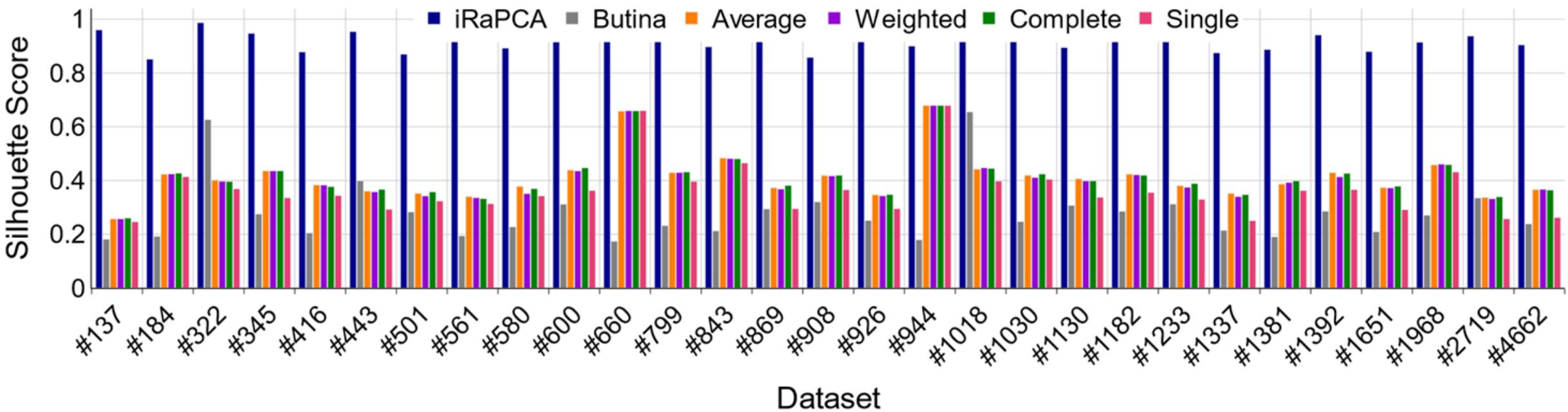
These plots are generated for each clustering round

Sunburst graph for the final clusters

# BENCHMARKING iRaPCA

Validation  
de iRaPCA

- Databases:
  - 29 from Cortes-Ciriano 2016
- Métrica:
  - Silhouette score





UNIVERSIDAD  
NACIONAL  
DE LA PLATA



**streamlit**  
**scikit-learn**  
**rdkit-pypi**  
**molvs**  
**mordred**  
**plotly**  
**seaborn**  
**validclust**  
**umap-learn**  
**openbabel**  
**scipy**  
**biopython**  
**pybiomed**  
**mols2grid**



Carolina Bellera



Lucas Alberca



Denis Prada



Manuel Llanos



Santiago Rodríguez



Maximiliano Fallico

# Muchas GRACIAS!



[www.lideb.biol.unlp.edu.ar](http://www.lideb.biol.unlp.edu.ar)



@lideb\_unlp



lideb\_unlp



<https://github.com/LIDeB>



questions?

