

Canviant la forma de veure els mètodes QSAR: aplicació de la projecció hiperbòlica a la química mèdica

Changing the way of viewing QSAR methods: the application of hyperbolic projection in medicinal chemistry

Roger Estrada, Santi Nonell i Jordi Teixidó
Universitat Ramon Llull. IQS School of Engineering

Resum: Els mètodes basats en lligands que permeten relacionar l'estructura molecular amb la predicció d'una propietat o activitat d'interès (QSAR/QSPR) són utilitzats en química mèdica per al disseny de fàrmacs amb unes propietats determinades. Els models obtinguts com a resultat acostumen a estar formats per un elevat nombre de variables que en dificulten la representació gràfica en l'espai euclidià. S'ha comprovat que el pas a un espai governat per una geometria hiperbòlica permet, mitjançant un canvi de mètrica, la representació d'espais multidimensionals sobre un espai de dues dimensions (en aquest cas, un disc de Poincaré). La projecció es realitza mitjançant l'optimització amb algorismes genètics de les coordenades de cada molècula en el disc, identificant aquelles que permeten mantenir la relació de semblança existent entre elles en l'espai original. Aquesta metodologia ha estat aplicada de forma complementària a altres mètodes de reducció dimensional per facilitar la selecció de les variables a emprar per a la predicció de la localització subcel·lular de fotosensibilitzadors.

Paraules clau: Geometria hiperbòlica, disc de Poincaré, espais multidimensionals, QSPR, QSAR, teràpia fotodinàmica.

Abstract: *Ligand-based methods used in medicinal chemistry permit mathematical relationships to be established between molecular structure and the prediction of a property or biological activity (QSPR/QSAR) in the design of drugs with specific characteristics. The resulting models are usually defined by a large number of variables, which complicates their graphical representation in Euclidean space. The use of hyperbolic geometry in the definition of space has been found to enable the representation of multidimensional points on a bidimensional space by changing the metric definition (in this case, using Poincaré's disk model). The projection of multidimensional points is performed by optimising genetic algorithms to find the coordinates which maintain intramolecular similarity. This methodology has been applied with other dimensional reduction techniques to identify the most appropriate variables to be used in the prediction of the subcellular localization of photosensitizers.*

Keywords: *Hyperbolic geometry, Poincaré's disk model, multidimensional spaces, QSPR, QSAR, photodynamic therapy.*

Breu introducció al disseny de fàrmacs

En els darrers anys, la descoberta de nous fàrmacs ha esdevingut un veritable repte per a la química mèdica. La incorporació de la química combinatoria dins de l'esquema sintètic, juntament amb els avenços realitzats en el camp del cribratge d'alt rendiment (HTS, *high-throughput screening*) per facilitar l'avaluació de l'activitat *in vitro* d'un gran nombre de compostos, va permetre que l'espai químic accessible per a la identificació de nous candidats a fàrmacs s'ampliés notablement.

Aquest fet comportà, emperò, dificultats associades al gran nombre de compostos accessibles, atès que no tots els punts de l'espai poden ésser conduïts al llarg de tot el procés de desenvolupament preclínic i clínic (principalment, a causa de motius econòmics, de temps i de recursos humans, entre d'altres).

Així, doncs, l'esquema tradicional utilitzat en el disseny de nous fàrmacs es caracteritza per la reducció progressiva de candidats. En les fases primerenques, una vegada es disposa d'una diana biològica validada, es proposa un conjunt de compostos prototip (*hit compounds*), dels quals s'espera que presentin una certa activitat biològica. Sobre aquest conjunt, que pot estar format per centenars de sistemes moleculars diferents, s'apliquen diverses tècniques (moltes d'elles, computacionals i/o estadístiques) per tal d'identificar-ne els millors o més representatius, per optimitzar-los i que puguin esdevenir compostos precandidats (*lead compounds*). Els poc més d'una dotzena de compostos precandidats que entren a les fases preclínica i clíniques són sotmesos a un seguit d'assajos

per demostrar la seva eficàcia en diferents models animals i, finalment, en humans. Idealment, després del cribratge de totes les etapes, tan sols un compost reïx i esdevé fàrmac.

Òbviament, l'explosió combinatòria de les etapes inicials on es desitja generar un espai químic el més divers possible es troba en contraposició de l'intent de discriminar els compostos indesitjables en les fases inicials (seguint el principi de *fail early*).¹ Tanmateix, l'espai químic disponible pot ésser definit, mitjançant tècniques computacionals, de forma virtual. Aquest fet posa en evidència la importància dels mètodes quimiinformàtics, que permeten predir la viabilitat dels compostos prototip de forma prèvia a la seva síntesi i a la consegüent avaluació de l'activitat biològica.

Actualment, les diverses tècniques computacionals que s'utilitzen per a aquesta finalitat es divideixen, segons el seu principal objecte d'estudi, en dos grans grups. Els mètodes basats en l'estructura utilitzen informació de la geometria tridimensional de la diana biològica per avaluar interaccions amb els lligands i poder identificar així aquells amb una major probabilitat de presentar l'activitat desitjada. Per la seva banda, els mètodes basats exclusivament en els lligands utilitzen la informació derivada de l'estructura molecular de diferents compostos actius i no actius, per tal de comparar-los i identificar les característiques comunes que els fan tenir l'activitat biològica desitjada.

Entre els últims mètodes esmentats, destaquen els mètodes QSAR/QSPR (*quantitative structure-activity/property relationships*), que tenen com a finalitat trobar una relació matemàtica entre algunes de les propietats dels lligands i la seva activitat biològica, per tal de poder-la predir de forma quantitativa. La informació fisicoquímica continguda en l'estructura dels sistemes moleculars és traduïda a valors numèrics, intel·ligibles pels sistemes informàtics, mitjançant representacions matemàtiques anomenades *descriptors moleculars*. El conjunt de tots els descriptors moleculars considerats constitueix les variables que permeten establir els models QSAR/QSPR de predicció i, consegüentment, definir numèricament l'espai químic.

Representació gràfica d'espais multidimensionals

La descripció matemàtica de l'espai químic es realitza normalment mitjançant un nombre superior a dos descriptors

que, juntament amb la funció de resposta, genera un espai n -dimensional que dificulta la seva representació i interpretació gràfica sobre el pla euclidià.

Actualment, es disposa de diferents mètodes que permeten reduir la dimensió d'aquest espai minimitzant la pèrdua d'informació inherent. En l'àmbit de la quimiinformàtica, és freqüent de recórrer a l'anàlisi de components principals (PCA, *principal components analysis*) per reduir els graus de llibertat dels models QSAR/QSPR a un conjunt de variables ortogonals (anomenades *components principals*, CP), definides com una combinació lineal de les originals² i que permeten explicar la major part de la variància d'aquestes (> 90 %, per exemple). La variabilitat de les dades conté, de fet, la informació que permet diferenciar els individus. Tot i així, sempre existeix la possibilitat que el nombre de CP necessàries per assolir el grau de variància explicada desitjada continuï essent major a tres, cosa que fa que l'espai químic no es pugui representar gràficament.

D'altra banda, les xarxes neuronals de Kohonen, anomenades també *self-organizing maps* (SOM), empren un entrenament no supervisat (per exemple, sense emprar informació referent al resultat esperat) per projectar espais multidimensionals sobre un pla format per neurones artificials. Tot i que conserven la topologia de l'espai inicial,³ la interpretació dels mapes obtinguts com a resultat no sempre esdevé una tasca senzilla.

Moltes de les dificultats atribuïbles a la interpretació dels mètodes de reducció dimensional es troben relacionades amb el fet d'estar acostumats a entendre la semblança molecular com la distància euclidiana entre dos punts de l'espai químic (determinat pels descriptors moleculars). Per aquest motiu, es planteja la possibilitat d'utilitzar eines geomètriques que permetin projectar les dades multidimensionals sobre un pla bidimensional i mantenir alhora les nocions de *distància*, *d'angle*, etc.

El disc de Poincaré

El fet que el postulat paral·lel d'Euclides no hagi estat demostrat com a teorema ha obert les portes a les anomenades *geometries no euclidianes*, entre elles, la geometria hiperbòlica.

Per a la representació gràfica d'aquestes geometries, hom pot recórrer a models de projecció de l'espai sobre el pla euclidià.

Un dels models més coneguts és el disc de Poincaré, l'ús del qual ja ha estat descrit per representar gràficament diversos tipus de xarxes multidimensionals.⁴⁻⁶

En aquest model, la totalitat del pla euclidià es representa dins d'un disc (H^2) de radi finit i unitari, definit sobre el camp dels nombres complexos:⁷

$$H^2 = \{z \in \mathbb{C} : |z| < 1\}$$

A diferència de la geometria euclidiana, des del punt de vista de la geometria hiperbòlica és possible, amb una recta r i un punt P determinats (figura 1), definir múltiples rectes paral·leles a r (és a dir, que no es creuen entre si) i que, alhora, passen per P .

Òbviament, aquesta situació comporta una modificació de la mètrica de l'espai, és a dir, de la manera com es mesura la distància entre dos punts. La mètrica de Poincaré estableix que la distància entre punts augmenta a mesura que hom s'allunya del centre. Aquest fet promou un efecte d'ull de peix (figura 2).

A més, la redefinició de la mètrica que governa l'espai comporta, a la vegada, la modificació de les trajectòries de mínima distància entre dos punts (és a dir, les rectes en l'espai

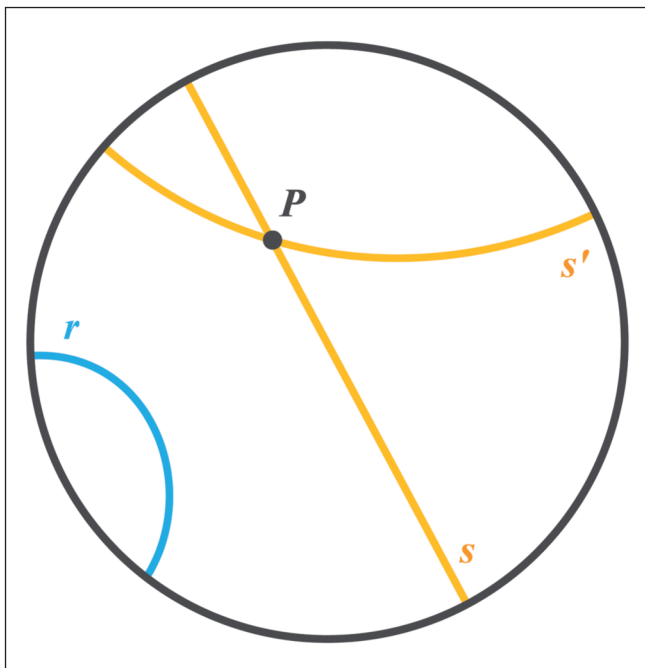


FIGURA 1. Representació esquemàtica de tres rectes (r , s i s') sobre el disc de Poincaré. Si bé les rectes s i s' s'intersequen en un punt P , cap d'elles no talla la recta r , de forma que (tant s com s') són alhora paral·leles a r .

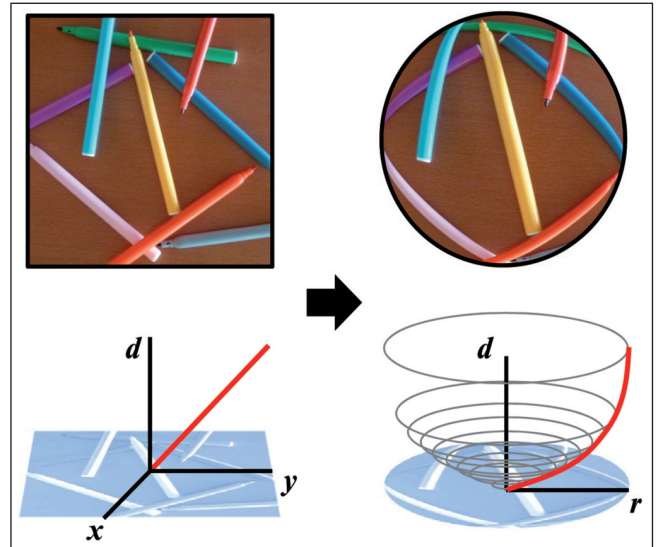


FIGURA 2. Representació gràfica del canvi de mètrica que té lloc en passar de la geometria euclidiana a la hiperbòlica. La distància entre qualsevol punt i l'origen de coordenades del pla euclidià augmenta de forma lineal segons el mòdul del vector posició (esquerra). En canvi, per encabir tot el pla euclidià dins d'un disc de radi unitari, la distància entre un punt del disc i el centre augmenta de forma no constant (dreta).

euclidià), així com dels operadors *translació* i *rotació* necessaris per manipular els punts dins del disc. En el pla hiperbòlic H^2 , aquestes operacions es realitzen per mitjà de les transformacions de Möbius⁸ (figura 3).

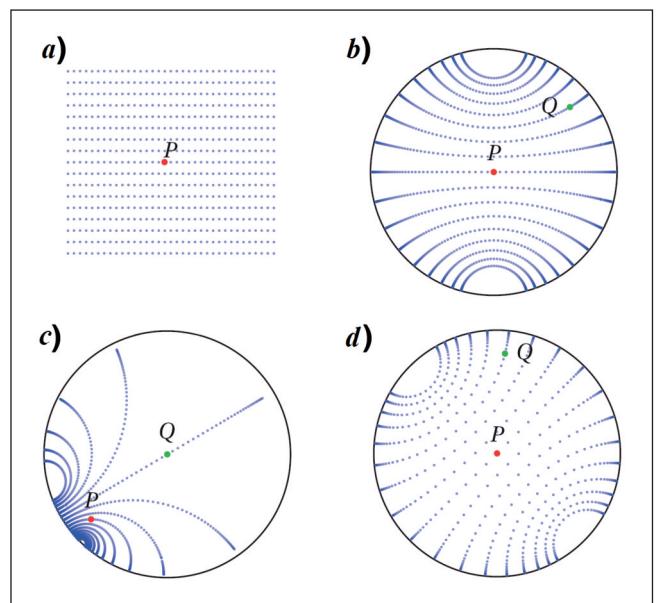


FIGURA 3. Amb un conjunt de rectes determinat, representades de forma puntejada, en el pla euclidià (a), el pas d'aquest sistema al disc de Poincaré provoca la curvatura de les seves trajectòries (b). Prenent com a referència els punts P i Q , les transformacions de Möbius permeten realitzar les operacions de translació de Q a l'origen (c) i de rotació del disc (d).

Optimització de la projecció sobre H^2

Amb l'objectiu de posar la informació anterior en el context de la quimioinformàtica, es considera un espai químic definit per n descriptors, generat per una determinada quimioteca virtual de compostos químics. Dins d'ell, cadascuna de les molècules pot ser entesa com un punt d'un espai n -dimensional $\{x_1, \dots, x_k\}$, que es pot representar sobre el pla H^2 . Per realitzar-ho, se cerquen aquelles coordenades complexes del disc de Poincaré $\{y_1, \dots, y_k\}$ on la distància hiperbòlica entre tots els punts implicats permeti mantenir la seva dissimilitud inicial.

Així, doncs, la projecció sobre H^2 requereix adoptar una definició de la semblança molecular. Com a químics, a partir de la representació estructural de les molècules orgàniques, hom pot tenir una certa intuïció de la similitud entre dues molècules pel que fa a la seva reactivitat, les seves propietats, etc. Tal com s'ha vist anteriorment, els descriptors són els encarregats de permetre als sistemes informàtics el fet de poder realitzar aquesta anàlisi segons el principi de semblança,⁹ és a dir, considerant que molècules semblants també ho seran pel que fa als descriptors i, per tant, presentaran propietats químiques similars. Això permet reduir la quantificació de la semblança molecular a la distància (normalment euclidiana, però no de forma exclusiva)¹⁰ entre els punts x_i .

Per tal de trobar les coordenades $\{y_1, \dots, y_k\}$ que millor s'ajusten en el disc de Poincaré, s'utilitzen algorismes d'optimització. Si bé els mètodes que utilitzen la informació del gradient per orientar la cerca del punt òptim (com, per exemple, el mètode *steepest descent*) han estat els utilitzats bibliogràficament per a aquesta finalitat,⁴ es proposa l'ús dels algorismes

genètics (com a mètode d'optimització estocàstic) per intentar millorar l'espai recobert durant la cerca.

Els algorismes genètics prenen suport en el concepte de l'evolució biològica per trobar, de forma iterativa, les millors solucions. Cadascuna de les possibles solucions (per exemple, les coordenades de la projecció) es considera un individu d'una població, en el qual els descriptors moleculars (o una codificació d'ells) esdevenen el material genètic que el defineix.

Com si es tractés realment d'un organisme biològic, els diferents individus d'una població generada aleatòriament es reproduïxen, recombinant el seu genoma, i donen lloc a nous organismes més ben adaptats al medi. Definint una funció (anomenada *fitness*) que permeti quantificar la bondat de l'adaptació de cada individu a l'entorn, la pressió selectiva exercida sobre una població inicial d'individus condueix a la supervivència de les millors solucions en les generacions successives. En el cas concret de la projecció sobre el pla H^2 , la funció de *fitness* correspon a la dissimilitud comentada anteriorment. Per a més informació sobre el funcionament i la implementació dels algorismes genètics, es recomana consultar el text de Goldberg.¹¹

S'ha pogut comprovar que l'aplicació dels algorismes genètics permet obtenir resultats del mateix ordre o superiors al mètode *steepest descent*. El mètode implementat s'ha validat amb diferents conjunts de dades per a les quals es disposava d'un alt coneixement previ o d'altres que es prenen normalment com a referència en tasques classificatòries, com és cas del conjunt de Fisher,^{12,13} basat en indicadors morfològics de tres espècies de flors d'iris.

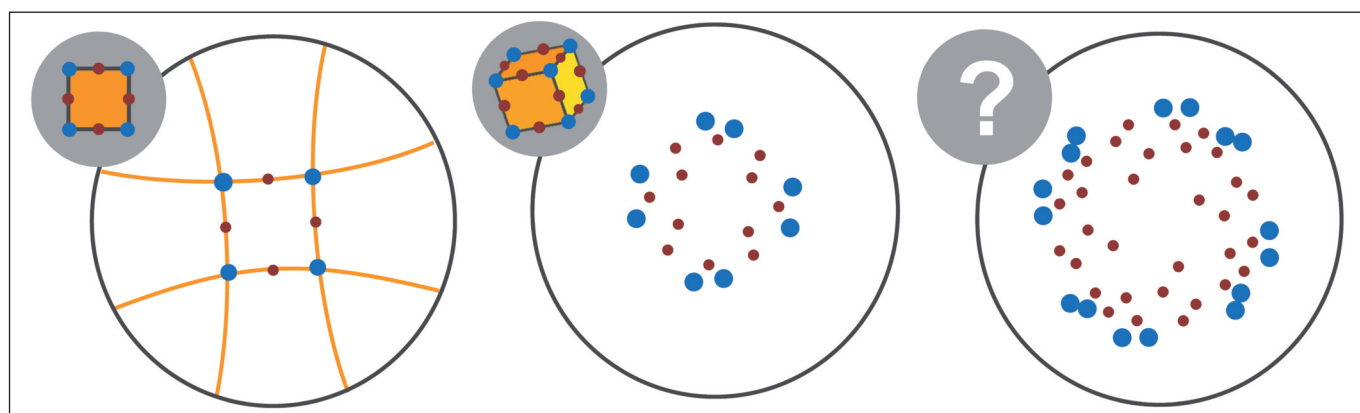


FIGURA 4. Aplicació de la projecció sobre el disc de Poincaré d'un quadrat (dues dimensions), un cub (tres dimensions) i un hiper-cub de quatre dimensions.

Atesa la independència que mostra la tècnica amb la naturalesa de les dades a projectar, aquesta pot ser emprada en un gran ventall d'aplicacions en les quals es requereixi projectar un gran nombre de dimensions sobre el pla, sempre que es disposi d'una definició de semblança vàlida. Pot utilitzar-se per a la representació de figures geomètriques complexes (com ara els hiperplans; vegeu la figura 4), en la visualització de resultats de *text mining*⁴ (i poder, per exemple, identificar fàcilment quina pel·lícula de ficció s'assembla més a quina altra, a partir de la freqüència de certes paraules del guió) o la semblança de les aigües embotellades comercials a partir de la seva composició química,¹⁴ entre d'altres. Tanmateix, en l'àmbit de la química mèdica, poden ser d'utilitat?

Aplicació de la projecció multidimensional en QSAR

La definició de quines variables (descriptors moleculars) són les més adequades per a l'establiment d'un model QSAR/QSPR de predicció és una de les tasques més difícils a l'hora d'aplicar aquests mètodes. Un nombre insuficient de variables pot conduir a una mala descripció del problema i fer que el model obtingut no sigui capaç d'establir correlacions entre les variables d'entrada i la propietat desitjada. D'altra banda, l'ús d'un nombre excessivament elevat de descriptors pot permetre una excel·lent descripció del sistema, que es tradueixi en una elevada capacitat de classificar correctament els elements del conjunt d'entrenament, però que, simultàniament, l'ajust sobre els elements sigui tan precís que el model sigui incapaç de realitzar prediccions quan se li presenta una entrada desconeguda¹⁵ (aquest comportament rep el nom d'*overfitting*).

En aquest sentit, la projecció dels descriptors utilitzats sobre el pla H^2 pot esdevenir d'una gran utilitat per aconseguir una interpretació gràfica del seu significat, atesa la conservació de la dissimilitud entre els punts.

Disseny de fotosensibilitzadors per a la teràpia fotodinàmica

La teràpia fotodinàmica (PDT, *photodynamic therapy*) aplicada al càncer es basa en la destrucció selectiva del teixit tumoral, aprofitant la fototoxicitat local derivada de la combinació de

llum visible, oxigen molecular i un fàrmac (anomenat *fotosensibilitzador*, FS).

Després d'administrar-se, el fàrmac es localitza preferentment en el teixit tumoral i, a continuació, la regió afectada és irradiada amb una font de llum de la longitud d'ona adequada. El FS és una molècula orgànica capaç d'absorbir aquesta llum (típicament, en la regió del vermell) i passar a un estat excitat de major energia. En tornar a l'estat fonamental, pot promoure el pas de l'oxigen molecular triplet, que hom té difós en tot l'organisme, a espècies reactives de l'oxigen (principalment, oxigen molecular en estat electrònic singlet). Aquestes espècies són capaces d'oxidar biomolècules de les cèl·lules malignes irradiades i provocar així la seva mort.

A causa de l'especificitat aconseguida per aquesta tècnica, nombrosos estudis bioquímics han demostrat la importància de la localització subcel·lular de l'agent FS en l'eficàcia global de la PDT.^{16,17} En aquests estudis, el mitocondri i els lisosomes han resultat ser els orgànuls cel·lulars que desencadenen d'una forma més eficient l'apoptosi cel·lular.¹⁸ Davant d'aquest fet, es planteja l'obtenció d'un model QSPR per a la predicció de la localització subcel·lular preferent de FS tetrapirròlics¹⁹ (figura 5).

A partir de la informació disponible bibliogràficament sobre la localització subcel·lular experimental d'aquests tipus de compostos, així com de models QSPR descrits prèviament

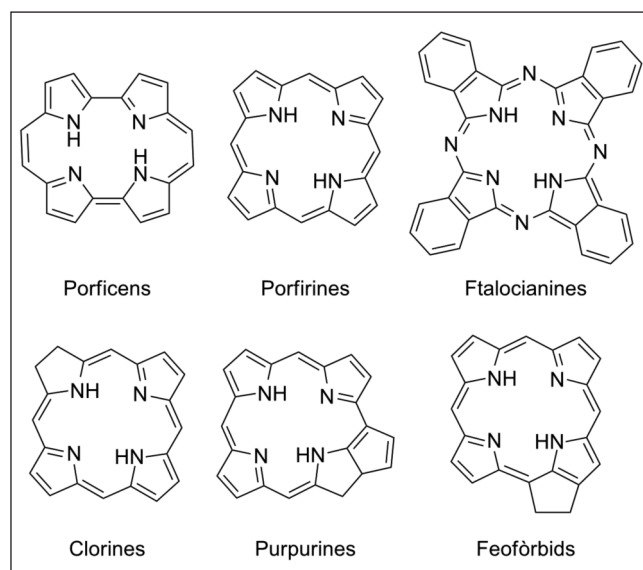


FIGURA 5. Fragments centrals dels fotosensibilitzadors tetrapirròlics de segona generació.

per a l'estudi de la selectivitat de fàrmacs en diferents òrgans,^{20,21} hom pot establir que la descripció matemàtica de la localització subcel·lular ha d'incloure informació referent al caràcter hidrofòbic dels FS, a la seva distribució de càrregues i a característiques estructurals relacionades amb la seva mida. Si bé aquestes característiques poden ser descrites per un gran nombre de descriptors moleculars (més d'un centenar), cal realitzar-ne una selecció per evitar la situació d'*overfitting*. Aplicant el mètode de selecció de variables *forward selection* (basat en la incorporació successiva de les variables) i prenent com a criteri de classificació l'algorisme *k-nearest neighbor*, s'aconsegueix reduir el nombre de variables a cinc. Per tal de poder discutir si aquests valors són capaços de descriure correctament les característiques estructurals de cadascuna de les famílies, se'n realitza la projecció sobre el pla H^2 (figura 6).

Tal com es pot observar, els descriptors seleccionats consideren químicament semblants clorines i feofòrbids, mentre que les diferencien de les ftalocianines (més apartades), tenint com a fil conductor les porfirines, que mostren una major diversitat estructural i que s'estenen al llarg de tot l'espai químic definit per aquests tres macrocicles (constitueixen, de fet, el nucli central dels altres tipus de compostos). Per la seva banda, els porfircens són els únics compostos que no contenen el nucli porfirínic a la seva estructura, i això fa que la seva

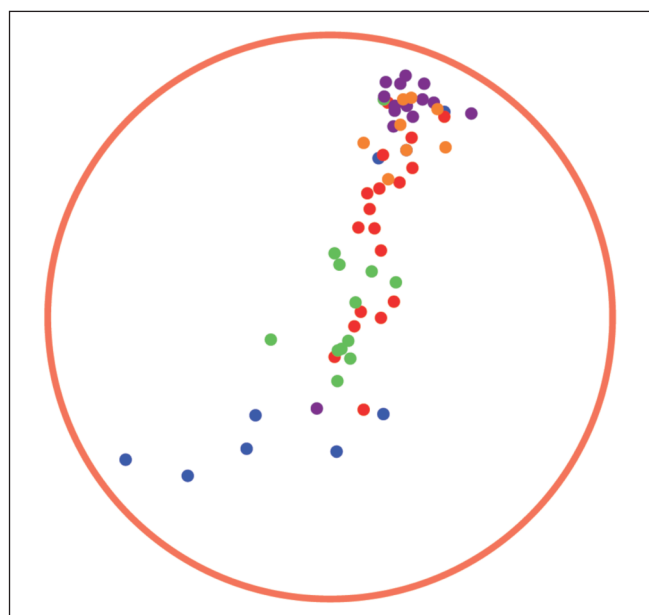


FIGURA 6. Projecció sobre el pla H^2 de l'espai de cinc descriptors emprats en el model QSPR de predicció de la localització subcel·lular de fotosensibilitzadors tetrapirròlics per a la PDT (● porfircens, ● porfirines, ● ftalocianines, ● clorines, ● feofòrbids).

projecció, en un dels extrems del diagrama, es trobi d'acord amb el sentit químic.

Si bé és cert que l'objectiu final del model QSPR és la predicció de la localització subcel·lular i no la classificació per famílies, la discussió dels resultats obtinguts en les projeccions mitjançant H^2 , SOM i PCA per a diferents grups de descriptors permet proposar un conjunt de descriptors crítics per a la correcta classificació subcel·lular. L'aplicació d'aquests descriptors per a l'entrenament d'un model basat en xarxes neuronals artificials permet obtenir un model QSPR de predicció de la localització subcel·lular amb una taxa d'encerts de prop del 85 %.

Conclusions

La projecció sobre un disc de Poincaré permet la representació gràfica de dades multidimensionals sobre un pla bidimensional, la qual cosa permet facilitar la interpretació de les dades originals de forma visual.

Aquesta estratègia, plantejada com l'alternativa als SOM, ha demostrat que és eficaç en la identificació de grups presents en dades multidimensionals, la qual cosa la fa idònia com a mètode de classificació o per avaluar l'adequació de les variables utilitzades per a la descripció de models QSAR/QSPR.

Agraïments

Aquest treball ha estat realitzat gràcies a la beca predoctoral concedida pel Comissionat per a Universitats i Recerca del Departament d'Innovació, Universitats i Empresa de la Generalitat de Catalunya i el Fons Social Europeu.

Referències

1. Boehm, M. *Virtual screening: Principles, challenges and practical guidelines*. Sottriffer, C. (ed). Wiley: Nova York, 2011, p. 3.
2. Gasteiger, J.; Engel, T. (ed.). *Chemoinformatics*. Wiley-VCH: Weinheim, 2003.
3. Kohonen, T. *Proceedings of the IEEE EICSTES 1990*, 78, 1464.

4. Walter, J.; Ritter, H. (ed.). *Proceedings of the Eighth ACM SIGKDD International Conference*. Edmonton. ACM: Nova York, 2002.
5. Bankston, D.; Battles, A.; Gurney, D.; Reyes, E. N.; Steidley, C. *Proceedings of the American Society for Engineering Education: Annual Conference and Exposition*. ASEE: Washington, 2003.
6. Ungar, A. A. *Proceedings of the 4th European SGI/CRAY MPP Workshop*. H. Lederer & F. Hertweck ed.: Garching, 1998.
7. Venema, G. A. *The Poincaré disk, colom 14*. Prentice Hall: Nova Jersey, 2005.
8. Girbau, J. *Geometria diferencial i relativitat*. UAB: Bellaterra, 1993. (Manuals de la Universitat Autònoma de Barcelona).
9. Johnson, M. A.; Maggiora, G. M. *Concepts and applications of molecular similarity*. Wiley: Nova York, 1990.
10. Pascual, R.; Borrell, J. I.; Teixidó, J. *Molecular Diversity* 2003, 6, 121.
11. Goldberg, D. E. «Genetic algorithms». A: *Search, Optimization and Machine Learning*. Addison-Wesley: Boston, 1989.
12. Fisher, R. A. *Annals of Eugenics* 1936, 7, 179.
13. Aydin, I.; Karakose, M.; Akin, E. *Applied Soft Computing* 2011, 11, 120.
14. Estrada, R.; Nonell, S.; Teixidó, J. VII Trobada de Joves Investigadors dels Països Catalans. Palma, 2012.
15. Tetko, I. V.; Livingstone, D. J.; Luik, A. *J. Chem. Comput. Sci.* 1995, 35, 826.
16. Stockert, J. C.; Cañete, M.; Juarranz, A.; Villanueva, A.; Horobin, R. W.; Borrell, J. I.; Teixidó, J.; Nonell, S. *Curr. Med. Chem.* 2007, 14, 997.
17. Richert, C.; Wessels, J. M.; Müller, M.; Kisters, M.; Benninghaus, T.; Goetz, A. E. *J. Med. Chem.* 1994, 37, 2797.
18. Nyman, E. S.; Hynninen, P. H. *J. Photochem. Photobiol. B* 2004, 73, 1.
19. Estrada, R.; Nonell, S.; Teixidó, J.; Sagristá, M. L.; Mora, M.; Villanueva, A.; Cañete, M.; Stockert, J. C. *Curr. Med. Chem.* 2012, 19, 2472.
20. Rashid, F.; Horobin, R. W.; Williams, M. A. *Histochem. J.* 1991, 23, 450.
21. Horobin, R. W.; Stockert, J. C.; Rashid-Doubell, F. *Histochem. Cell Biol.* 2006, 126, 165.



R. Estrada



S. Nonell



J. Teixidó

Roger Estrada és llicenciat en química i enginyer químic de l'Institut Químic de Sarrià (Universitat Ramon Llull, URL). Doctor per la URL, ha participat en diversos projectes de recerca de química mèdica (disseny molecular i estudis de química teòrica per a la recerca de nous fàrmacs).

Santi Nonell és professor catedràtic del Departament de Química Orgànica de l'Institut Químic de Sarrià (Universitat Ramon Llull) i responsable del laboratori de fotoquímica del Grup d'Enginyeria Molecular.

Jordi Teixidó és professor catedràtic del Departament de Química Orgànica de l'Institut Químic de Sarrià (Universitat Ramon Llull) i responsable del laboratori de disseny molecular del Grup d'Enginyeria Molecular.