

Representation in Machine Learning for Drug Design

**Ross D. King,
University of Manchester & Alan Turing Institute
ross.king@manchester.ac.uk**

Data Representation

Effective Data Representations

- n The key to success in machine learning is the use of effective data representations.
- n Almost all machine learning is based on representations that use tuples of attributes, i.e. the data can be put into a single table, with the examples as rows, and the attributes (descriptors) as columns.
- n Attributes are normally intrinsic properties of the examples believed to be important: for example if one wished to learn about the effectiveness of a drug, then properties of its molecular structure may be useful attributes.

Notes 1 - Deep Learning

- n The most exciting current area of machine learning is that of deep neural networks (DNNs).
- n The success of DNNs has been based on their ability to utilize multiple neural network layers, and large amounts of data, to learn how to convert raw input representations (e.g., image pixel values) into richer internal representations that are effective for learning.
- n Thanks to this ability to learn effective internal representations, DNNs have succeeded in domains that had previously proved recalcitrant to ML.

Notes 2 – Relational Learning

- n The use in machine learning of representations based of tuples of attributes is essence the use of propositional logic.
- n Use of the richer, more expressive, language of 1st-order predicate logic is termed Relational Learning.
- n Not fashionable now, but in my view, sooner or later relational learning will become essential.
- n Used to be considered too inefficient to be practical, but DNNs have moved the bar here!

Talk Plan

- n Transformative Machine Learning.
- n Relational Machine Learning

Transformative Machine Learning

**Ross D. King,
University of Manchester & Alan Turing Institute
ross.king@manchester.ac.uk**

**Multiple Related Problems
are Typical in Science**

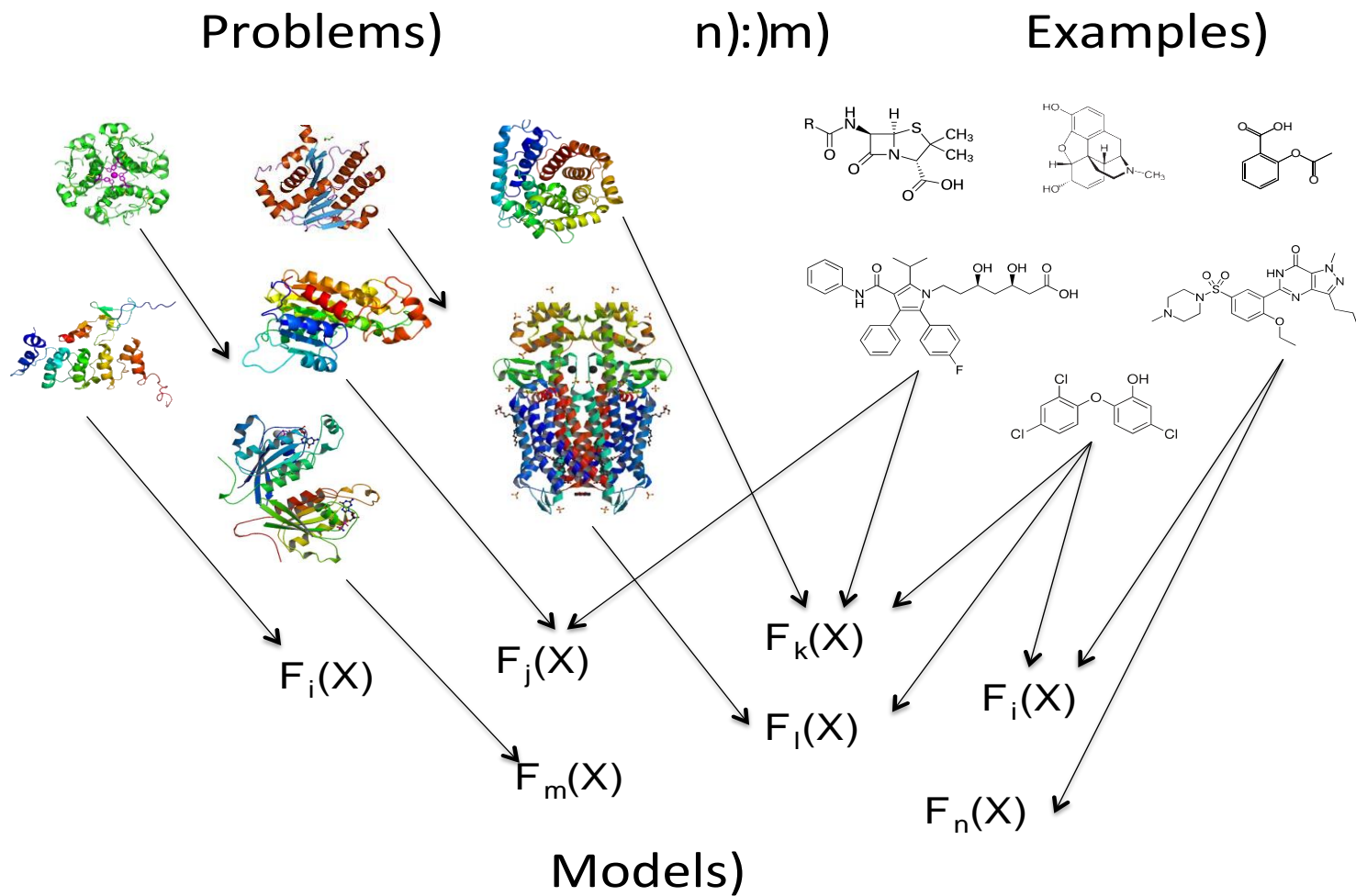
Multi-task learning

- n Scientific problems often present themselves as sets of related problems
- n Multi-task learning is the branch of machine learning in which related problems (tasks) are learned simultaneously.
- n The the aim is to exploiting similarities between the problems to improve performance.
- n Problems are learned in parallel using a shared representation, so what is learned from one task can also be used for another problem.

Transfer Learning

- n Transfer Learning is closely related to Multi-task learning.
- n Information is transferred from a specific source problem to a specific target problem.
- n This can be achieved by forcing the target model to be structurally or otherwise similar to the source model(s).
- n The success or failure of multi-task or transfer learning often crucially depends on the existence of a good task similarity measure.

Prediction Problems

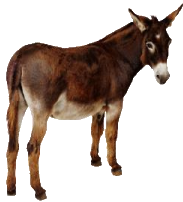


Transformative Learning

The New Idea

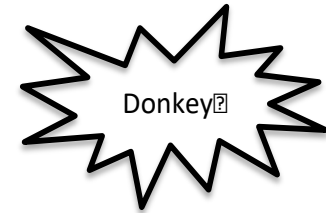
- n The fundamental new idea is to transform the representation from one based on describing examples using intrinsic properties, to an extrinsic representation based on what other models predict about examples.
- n This transformation has the dual advantages of: producing significantly more accurate predictions, and providing explainable models.

Standard Machine Learning



+ **?**

Size?	Ears?	Cute?	Donkey?
Big?	Big?	No?	1.0?
Small?	Big?	No?	0.3?



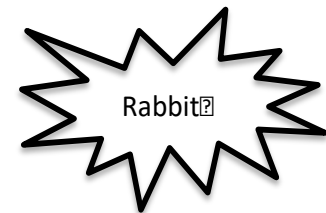
+ **?**

Size?	Ears?	Cute?	Kitten?
Small?	Small?	Yes?	1.0?
Small?	Small?	No?	0.1?

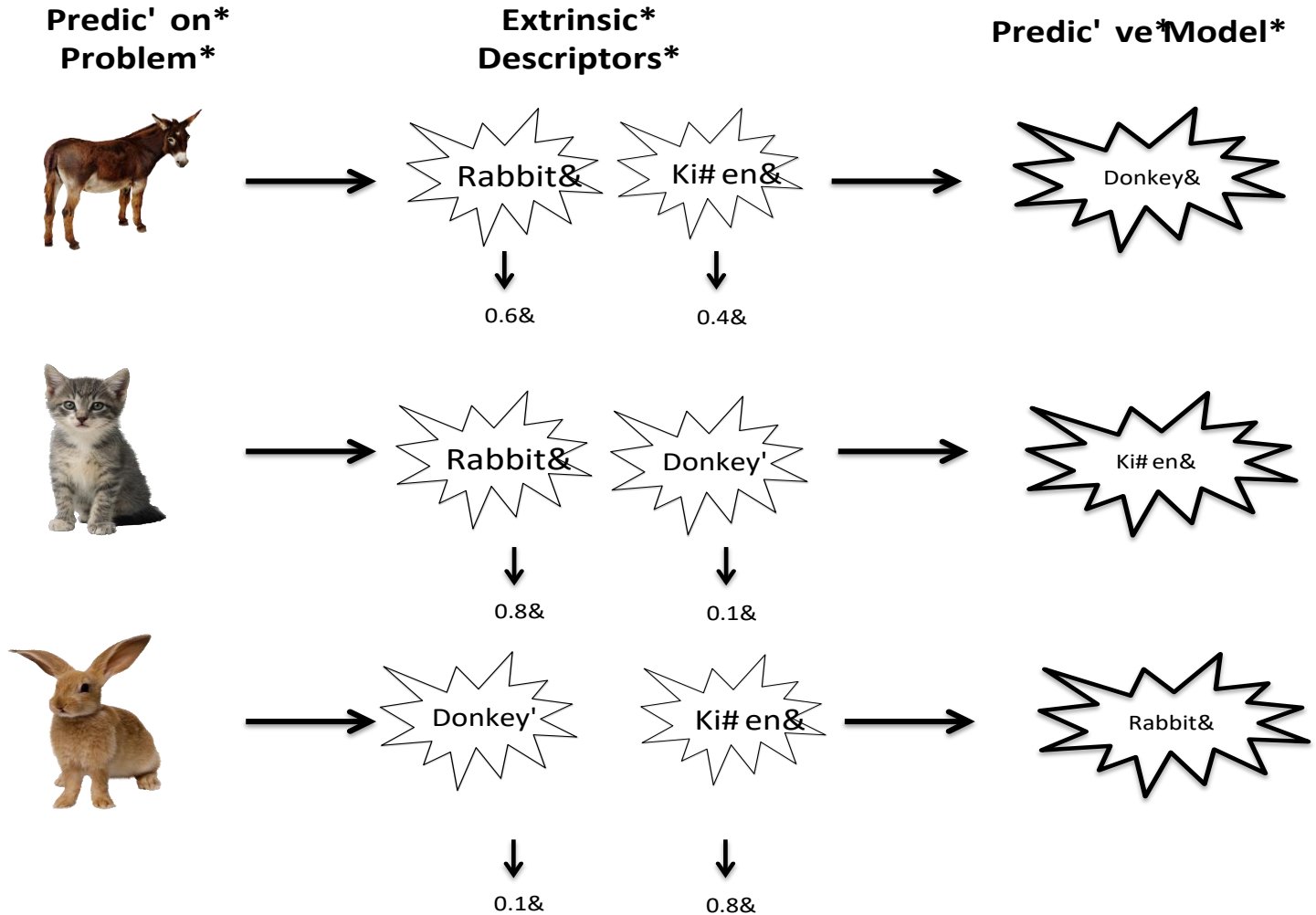


+ **?**

Size?	Ears?	Cute?	Rabbit?
Small?	Big?	Yes?	1.0?
Big?	Small?	Yes?	0.2?



Transformative Machine Learning

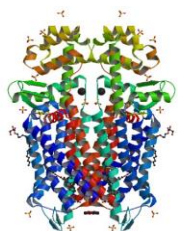


Standard Machine Learning

Predic' on*Problem*

Intrinsic*Descriptors*

Predic' ve*Model*



+ "

	G1#	G2#	...#	G1024#	Activity#
<chem>Oc1ccc(O)c2c1O[C@H]3[C@@H](C)N[C@@H](C)C3</chem>	#	#	#	#	#
	1#	0#	1#	1#	0.9#
...#	...#	...#	...#	...#	...#
<chem>CC(C)S[C@@H]1[C@H](C(=O)N1C(=O)O)C(=O)R</chem>	#	#	#	#	#
	0#	0#	1#	1#	0.3#

→

$F_1(X)$

"

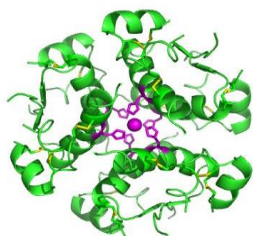
...

"

...

"

...



+ "

	G1#	G2#	...#	G1024#	Activity#
<chem>Oc1ccc(O)c2c1O[C@H]3[C@@H](C)N[C@@H](C)C3</chem>	#	#	#	#	#
	0#	0#	1#	0#	0.5#
...#	...#	...#	...#	...#	...#
<chem>CC(C)S[C@@H]1[C@H](C(=O)N1C(=O)O)C(=O)R</chem>	#	#	#	#	#
	1#	0#	1#	1#	0.4#

→

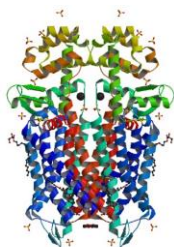
$F_n(X)$

Transformative Machine Learning

Prediction Problem

Extrinsic Descriptors

Transformative Learning
Prediction Model



+

	\mathcal{F}_1	\mathcal{F}_2	...	\mathcal{F}_n	Activity
	0.3	0.4	0.7	0.9	0.9
...
	0.1	0.8	0.7	0.5	0.3

→

$T_1(X)$

||

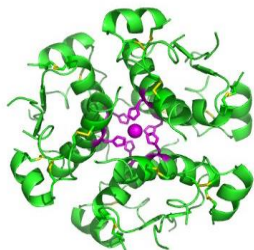
||

||

...

...

...



+

	\mathcal{F}_1	\mathcal{F}_2	...	\mathcal{F}_n	Activity
	0.2	0.5	0.9	0.1	0.5
...
	0.5	0.8	0.9	0.1	0.4

→

$T_n(X)$

Intuition - 1

- n The intuition behind Transformative Learning is that instead of describing examples by intrinsic properties, examples are described by what other models predict about them.
- n This combining ideas from multi-task learning and transfer learning.

Intuition - 2

- n Instead of using a predefined similarity measure to pre-select a set of similar tasks, we project the different tasks into one joint numeric representation.
- n Then use a meta-learning algorithm to learn from this new representation how to make accurate predictions for the task at hand.
- n Transfer learning has the advantage that it is machine learning method agnostic, and can be applied to improve the performance of many different methods.

The Standard Alternative

- n The standard alternative multi-task machine learning approach is to try to learn one large model that encompasses all the problems.
- n In many circumstances this can work well.
- n However, this standard approach has the disadvantages that if new data, or a new problem is added, then the whole model needs to be relearned.
- n Also, neither the relationships between problems, nor the relationships between examples, is made explicit.

Methods

Experimental Setup – Machine Learning Methods

- n To investigate the utility of Transformative Learning we selected four machine learning methods:
 - Random forests (RF).
 - Support-vector machines (SVMs).
 - k-nearest neighbour (KNN).
 - Neural-networks (NN).
- n These represent the main families of non-linear Machine Learning methods.

Experimental Setup – Problems

- n We applied the four machine learning methods to three typical real-world scientific problems:
 - Drug-design (quantitative structure activity relationship learning - QSAR).
 - Predicting human gene expression (across different tissue types and drug treatments) – LINCS.
 - Meta-machine learning (predicting how well machine learning method will work on problems).

Results

Results

Problem Area	ML Method	Standard NRMSE (RMSE)	Transformed NRMSE (RMSE)	Difference % NRMSE (RMSE)
<i>QSAR</i>				
	<i>RF</i>	0.1643 (0.6511)	0.1478 (0.6316)	10.05 (3.07)
	<i>SVM</i>	0.1693 (0.6720)	0.1522 (0.6490)	10.10 (3.42)
	<i>KNN</i>	0.167 (0.711)	0.171 (0.734)	-2.39 (-3.23)
<i>LINCS</i>				
	<i>RF</i>	(0.0694)	(0.0664)	3.68 (4.32)
	<i>SVM</i>	(0.0692)	(0.0677)	10.9 (2.18)
	<i>KNN</i>			2.92 (-0.28)
	<i>NN</i>	(0.0742)	(0.0707)	-1.15 (4.72)
<i>Meta-ML</i>				
	<i>RF</i>	0.257 (0.1184)	0.124 (0.0526)	51.8 (55.57)
	<i>SVM</i>	0.296 (0.1340)	0.214 (0.0972)	27.8 (27.2)
	<i>KNN</i>	0.290 (0.1330)	0.274 (0.1260)	5.39 (5.40)
	<i>NN</i>	0.323 (0.1480)	0.248 (0.1130)	23.5 (23.7)

XAI - Transformative Learning - 1

- n Explainable AI (XAI) is an increasingly important area of machine learning, for in many applications (e.g. medical, financial) there is a necessity to make predictions understandable.
- n The understandability of ML models depends on model simplicity, and on how closely the model reflects human concepts.
- n The standard theory of human concepts dates back to Aristotle, and is based on the presence of necessary and sufficient conditions.

XAI - Transformative Learning - 2

- n The understandability of Transformative learning models is based around the cognitively natural concepts of problem similarity and prototypes.
- n This approach to describing human concepts has been a focus for much recent cognitive science research.
- n The Transformative Learning representation is particularly useful if the intrinsic attributes are unsuitable for understandable modelling.

XAI –Human DHFR

	Target ID	Name	Species
1	CHEBL2902	Dihydrofolate reductase	<i>Lactobacillus casei</i>
2	CHEBL2014	Nociceptin receptor	<i>Homo sapiens</i>
3	CHEMBL2111414	Tyrosine-protein kinase ABL	<i>Homo sapiens</i>
4	CHEMBL3048	Nitric-oxide synthase, brain	<i>Rattus norvegicus</i>
5	CHEMBL329	Type-1 A angiotensin II receptor	<i>Rattus norvegicus</i>
6	CHEMBL4264	24-sterol C-methyltransferase	<i>Arabidopsis thaliana</i>
7	CHEMBL5457	Dihydrofolate reductase	<i>Mycobacterium avium</i>
8	CHEMBL5372	Methionyl-tRNA synthetase	<i>Staphylococcus aureus</i>
9	CHEMBL1075294	Indoleamine 2,3-dioxygenase 1	<i>Mus musculus</i>
10	CHEMBL1741172	60S ribosomal protein L 19-A	<i>Saccharomyces cerevisiae</i>

XAI – Drug Design Example- 1

- n Standard – Explain activity in terms of molecular features.
- n Transfer learning – Explain activity in terms of similarity of related problems.
- n As expected other DHFR target models are recognised as being useful in predicting human DHFR activity.

XAI – Drug Design Example- 2

- n The other target models selected are of biological interest.
- n It is interesting that *H. sapiens* Tyrosine-protein kinase ABL was selected, as it has been empirically shown that it is possible to jointly target *H. sapiens* DHFR and tyrosine kinases.
- n It is also intriguing that Tetrahydrobiopterin (BH4) is a required cofactor for nitric oxide synthase, and that DHFR regenerates 7,8-dihydrobiopterin (BH2) into BH4.

Large-scale Modelling

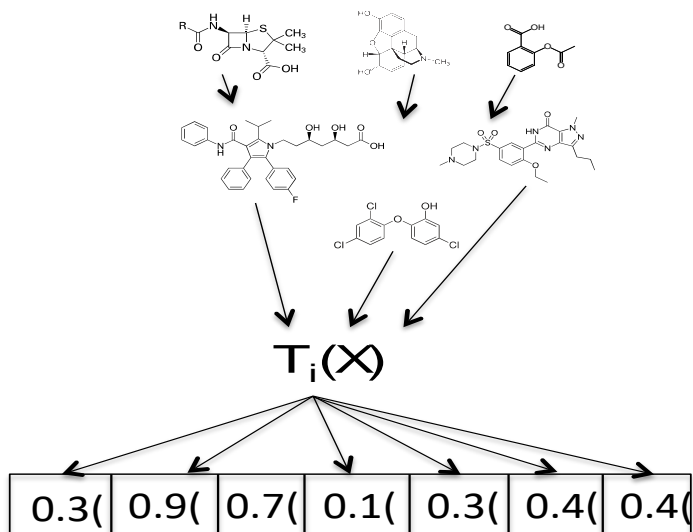
- n An important side-product of TL is the large-scale production of predictive models.
- n These can be used for clustering, and to make large-scale predictions.

Gene/Protein Clustering

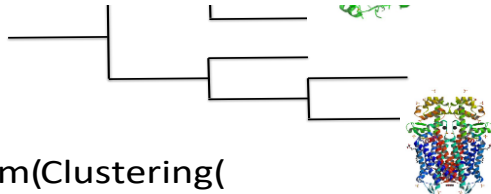
- n The standard approach to estimating gene/protein similarity is to estimate evolutionary distance by sequence comparison.
- n What is most important in most practical problems is not evolutionary distance, but functional similarity.
- n Use of the Transformative Learning predictive models enable genes/proteins to be compared using the accumulated information from millions of empirical experiments.

Clustering

Problem(Profile(

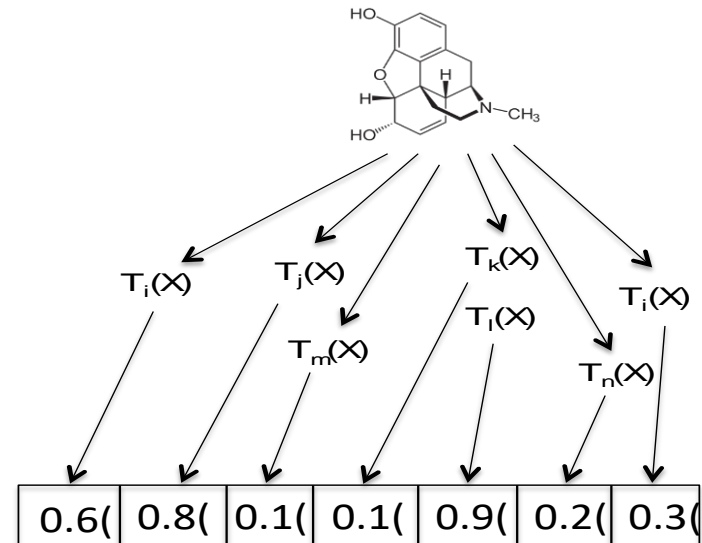


1(Problem((X((n(Examples(

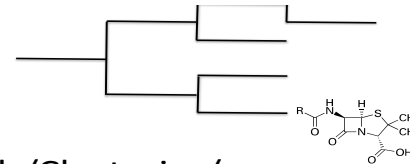


Problem(Clustering(

Example(Profile(



1(Example((X((n(Problems(



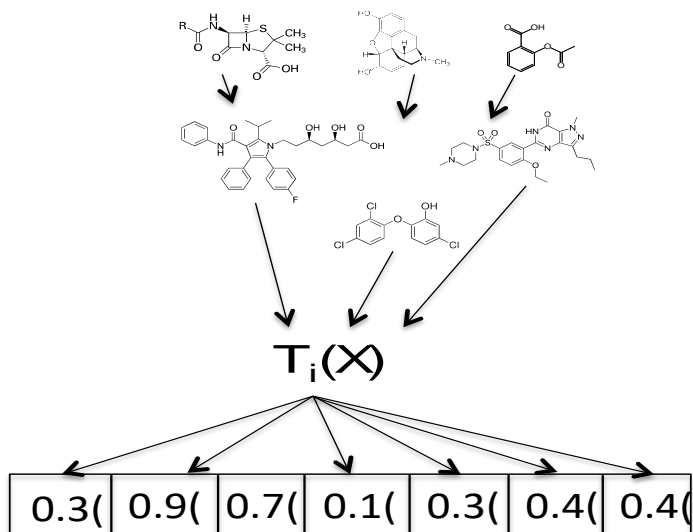
Example(Clustering(

Compound Clustering

- n A fundamental problem in chemoinformatics is the estimation of the similarity of chemical compounds.
- n Many different approaches have been applied based on chemical structure.
 - The most commonly used such methods are Tanimoto (Jaccard) coefficient distance between molecular fingerprints, and graph similarity.
- n What is of central interest in most chemoinformatic applications is not structural similarity, but functional similarity.

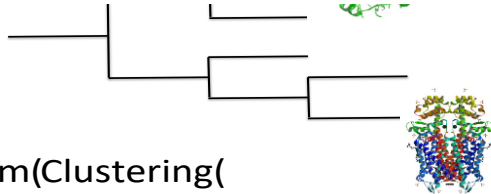
Clustering

Problem(Profile(

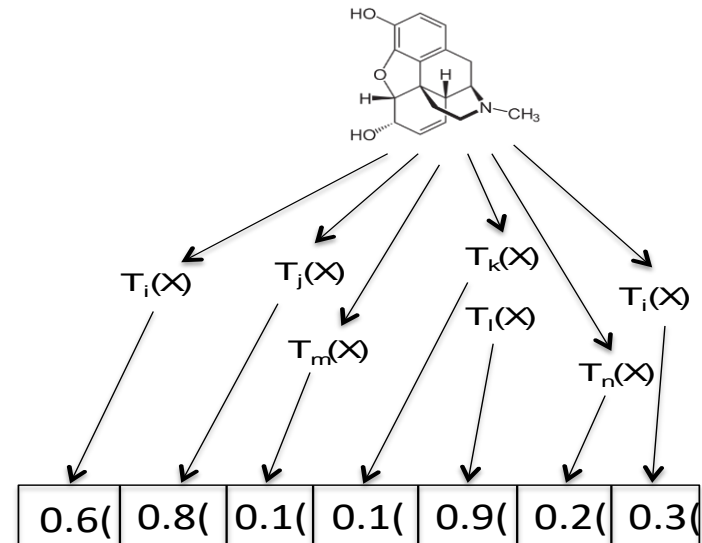


1(Problem((X((n(Examples(

Problem(Clustering(

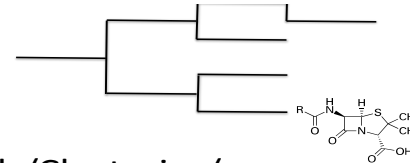


Example(Profile(



1(Example((X((n(Problems(

Example(Clustering(

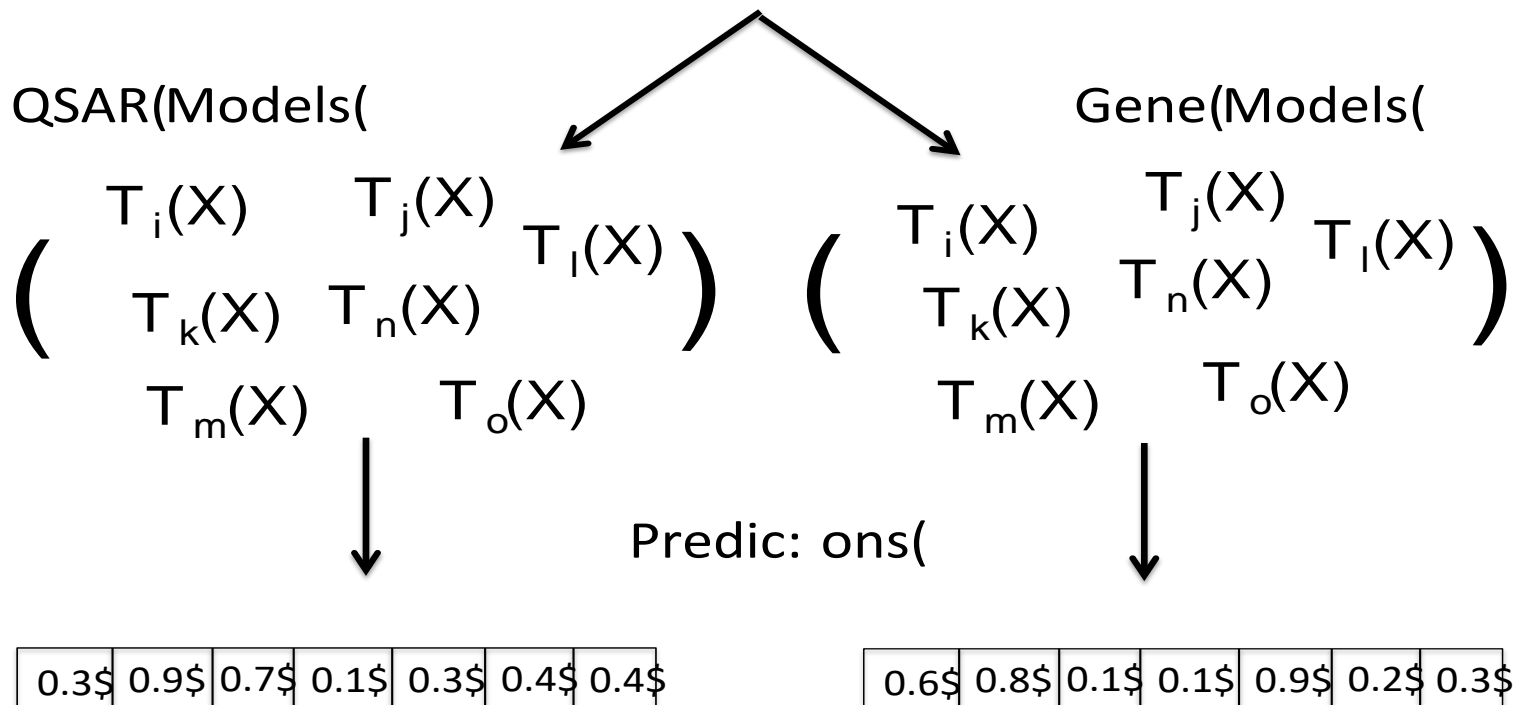


Large-scale Predictions

- n The multiple prediction models of Transformative Learning enable large-scale model predictions.
- n To illustrate this we have applied our thousands of QSAR target models to predict all the compounds found in ChEMBL
- n The universe of possible small-molecules is many many orders of magnitude larger than ChEMBL compounds.
- n We therefore plan to bundle up the prediction code and made it freely available.

Large-scale Predictions

{ Universe (of small molecules) }



FAIR Data

- n All the models, predictions, and clustering are fully annotated with meta-data, linked, and freely published on the semantic web.
- n To maximize its added-value we follow the FAIR (Findability, Accessibility, Interoperability, and Reusability) principles for publishing digital objects.
- n The goal is to facilitate discovery, evaluation, data and knowledge integration and reuse by the community after the publication process.

Discussion

Extensions - 1

- n I have presented only the simplest form of Transformative learning. Many other variants are possible with likely improved performance.
- n One such approach would be to apply feature selection. The lack of feature selection may explain the K-nn results,

Extensions - 2

- n Ensemble Transformative Learning - apply multiple ML methods to obtain better performance than could be achieved by any single ML method alone.
- n One possible approach would be to include the standard intrinsic descriptors as well as the Transformative Learning descriptors in the Transformative learning step.
- n Stacking could also be used to learn the best ensemble Transformative learning method from different feature selection and ensemble methods.

2nd Order Transformative Learning

- n It is natural to extend the idea of Transformative Learning by applying it a second time
 - i.e. to use the predictions from the transformed representation to form a second-order transformed representation.
- n As the predictions from the transformed representation are better than the ones from intrinsic representation, learning using second-order transformed representation can be more successful than with the first -order transformed representation.

Conclusions

- n We have presented a novel and general Machine Learning representation for sets of related problem.
- n The representation has the dual advantages of improving performance, and enabling explainable predictions.
- n The fundamental new idea is to transform a standard data representation into an explicit representation based on the predictions of pre-trained models.
- n As Machine Learning is increasing being applied to large sets of related problems, we expect Transformative Learning to be of broad general application to scientific problems, and beyond.

Acknowledgments

- n Ivan Olier, Oghenejokpeme I. Orhobor, Larisa N. Soldatova, Joaquin Vanschoren.
- n This research was partly supported by the Engineering and Physical Sciences Research Council (EPSRC) grant EP/K030469/1.

Acknowledgments

- n Ivan Olier, Oghenejokpeme I. Orhobor, Larisa N. Soldatova, Joaquin Vanschoren.
- n This research was partly supported by the Engineering and Physical Sciences Research Council (EPSRC) grant EP/K030469/1.

Relational Machine Learning

**Ross D. King,
University of Manchester & Alan Turing Institute
ross.king@manchester.ac.uk**

Attributes

- n Most statistical, neural network, machine learning and data mining methods use attributes to represent examples.
- n An attribute is something true about the whole example.
- n Example can be put into a single row in a Table.

Attributes and chemical structure

- n Difficult to represent arbitrary chemical structure using a single table. Either the table grows exponentially, or a compromise must be made in fidelity.
- n What should the attributes be?
 - Hansch
 - Boolean fingerprints
 - topological indices
 - voxols

Relational Learning

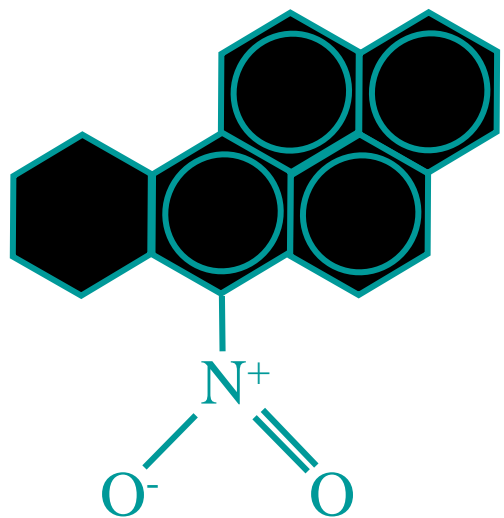
- n Uses richer representation language
 - allows relations and therefore structure
 - based on computer programs (predicate logic)
 - no need to use only a single row of a table
 - easily extended to 3D structure, multiple conformations, etc.

Mutagenesis Example

Example: Predicting Mutagenesis

- n Structure-activity relationships of mutagenic aromatic and heteraromatic nitro compounds.
- n Aromatic nitro compounds have been found to be both mutagenic and carcinogenic.
- n The mutagenesis of these compounds has been studied using the Ames test.

Example Compounds



6-nitro-7,8,9,10-tetrahydro[a]pyrene



Hansch Approach 1

- n Corwin Hansch was one of the pioneers of QSAR.
- n He developed the classical method for QSAR of using linear regression based on describing compounds using attributes such as logP and indicator variables.

Hansch Approach 2

- n After much experimentation on the mutagenesis data the following equation was proposed by Hansch and co-workers.
- n
$$\text{Log TA98} = 0.65(\pm 0.16)\log P - 2.90(\pm 0.59)\log (b.10\log P + 1) - 1.38(\pm 0.25)\text{Lumo} + 1.88(\pm 0.39)|1 - 2.89(\pm 0.81)|a - 4.15(\pm 0.58)$$
- n $n=188, r=0.900, s=0.866, \log PO=4.93, \log b=-5.48, F_{1,181}=48.6$
- n This equation gives little structural insight into the QSAR.

Relational Approach

- n The compounds are represented using just the atom and bond information.

Relational Approach: Atoms

- n The predicate `atomm` is used, e.g. in compound 127 (3,4,3'-trinitrobiphenyl).
- n `atom(127, 127_1, c, 22, 0.191)`.
- n States that in compound 127, atom no. 1 is a carbon of type 22 (aromatic carbon in a six membered ring) with a partial charge of 0.191.

Relational Approach: Bonds

- n The predicate bond is used, e.g.
bond(127, 127_1, 127_6, 7).
- n States that in compound 127, atom no. 1 and atom no. 6 are connected by a bond of type 7 (aromatic bond).
- n There are ~18,300 atom and bond facts in the background knowledge.

High Level Background

Knowledge 1

- n Relational Learning enables the inclusion of a priori background knowledge.
- n The background knowledge consists of computer programs. Arbitrary complex pieces of prior knowledge can be used in the learning process.
- n Generally the input programs are in the language PROLOG; but in principle could be any language that can be linked to PROLOG, e.g. FORTRAN.
- n It is like using an expert system in the learning process.

High Level Background

Knowledge 2

A PROLOG definition of a Methyl group is:

methyl(Drug,[Atm0,Atm1,Atm2,Atm3,Atm4]) :

atom(Drug, Atm0, Type,.....),	% Link atom
Type not h,	
atom(Drug, Atm1, c. 10, ...),	% Aryl carbon
atom(Drug, Atm2, h, 3, ...),	% Hydrogen
atom(Drug, Atm3, h, 3, ...),	% Hydrogen
atm(Drug, Atm4, h, 3, ...),	% Hydrogen
bond(Drug, Atm0, Atm1,1),	% Single bond
bond(Drug, Atm1, Atm2, 1),	% Single bond
bond(Drug, Atm1, Atm3, 1),	% Single bond
bond(Drug, Atm1, Atm4, 1),	% Single bond
Atm2 @> Atm3,	
Atm3@> Atm4.	% Atm2,3,4 different

High Level Background

Knowledge 3

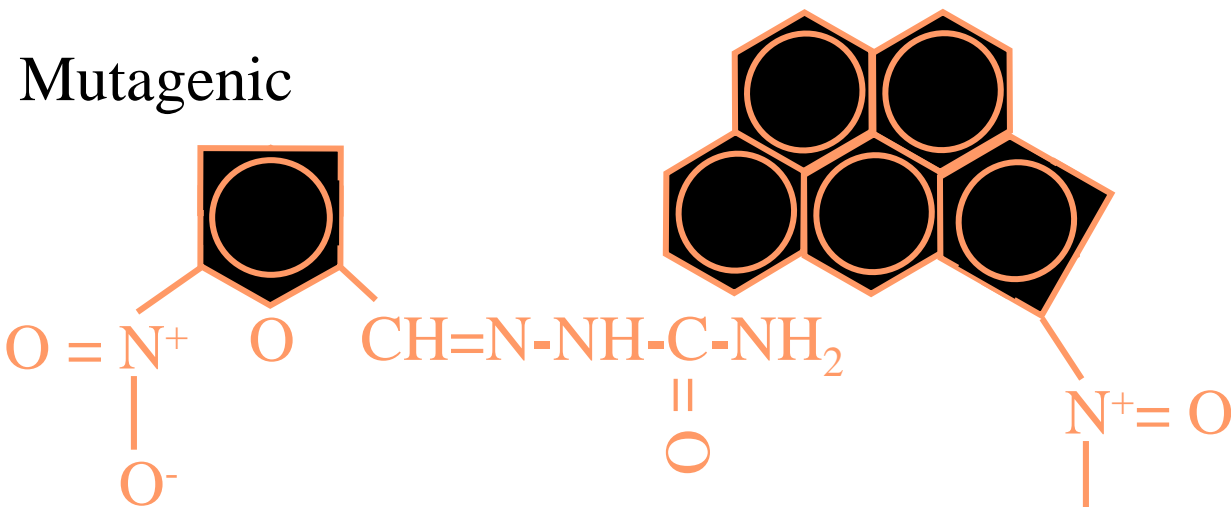
- n It is important to appreciate that encoding Progol programs to define concepts is not the same as including them as 'indicator variables'.
- n This method is fully automatic.
- n Relational learning can use structural combinations of the chemical groups, e.g. could in theory learn that a structural indicator for activity is diphenylmethane (a benzene single bonded to a carbon atom single bonded to another benzene).

Results 1

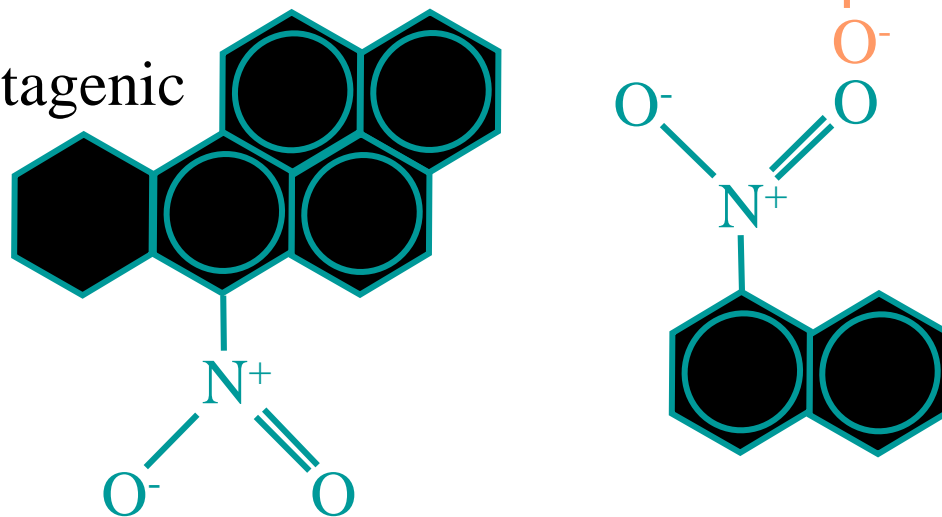
- n Hansch split the data into two subsets: 188 compounds which could be explained by a regression equation, and 42 compounds which could not.
- n For the 188 compounds a SAR was found that was as good as Hansch.
- n For the 42 compounds a SAR was found that was significantly ($P < 0.025$) more accurate than linear regression, quadratic regression, or neural nets.

ILP Structural Alert for Mutagenicity

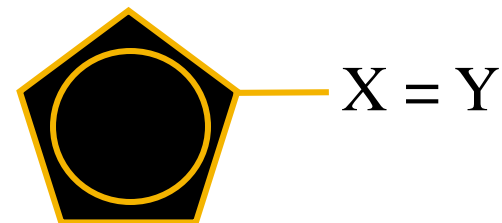
Mutagenic



Non-Mutagenic



Discriminating
Pattern



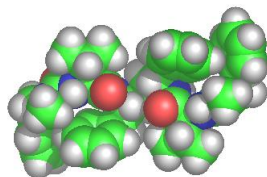
Extensions

- n Structure-Based Drug Design
 - Use the 3-dimensional structure of the drug (ligand) and protein.
 - Multiple conformations.
- n Quantum Mechanics
 - Molecules are quantum mechanical objects, not classical ones. We wish to represent this quantum mechanical while learning.
- n Data Mining
 - Find frequent patterns in sets of compounds.

Three Dimensions

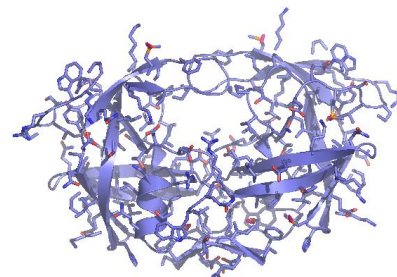
Structure-Based Drug Design

Ligand-based design



Ligands with known activity profile

Structure-based design



Target of known 3D structure

- Features involved in the binding mode to develop new compounds
- Most likely orientation of the ligand in the target (*docking*)
- Estimation of the strength of interaction (*scoring*)

Glycogen phosphorylase *b*

A case study : learning from a series of 3D complexes

Application to 51 structures of the glycogen phosphorylase *b*

n **Water**

Major problem for most of the 3D QSAR techniques

```
:- modeb(*, water(+drug, -conf, -pos)).
```

n **Amino acids**

Computational cost

✓ Properties

✓ Defined active site

```
:- modeb(*, prot_hacc(+drug, -conf, #aa, -pos)).  
:- modeb(*, prot_backc2(+drug, -conf, #aa, -pos)).  
:- modeb(*, prot_amide(+drug, -conf, #aa, -pos)).  
:- modeb(*, prot_lipo_seg(+drug, -conf, #aa, -pos)).  
...
```

n **Specific interactions**

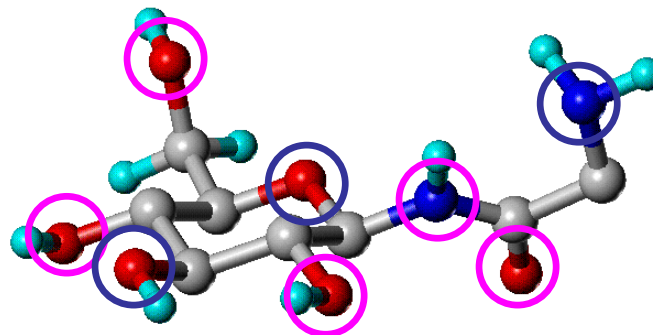
```
:- modeb(*, hb(+drug, -conf, -pos, -pos)).  
:- modeb(*, ionic(+drug, -conf, -pos, -pos)).
```

Ligand-based study

51 inhibitors in their crystallized conformation : ideal situation for 3D-QSAR

	R^2_{cv}
CoMFA	0.46
ILP	0.66
ILP/CoMFA	0.54

- ILP > ILP/CoMFA > CoMFA
- Most functional groups involved in the model
- Amide group in that position is related to highly active
- P3 negative contribution



$$\log(1/K_i) = 2.41 + 0.84 \cdot P1 + 0.95 \cdot P2 - 0.40 \cdot P3$$

P1 : hdon (C) , hacc (D) , hacc (E) , amide (F) ...

P2 : hdon (C) , hacc (D) , hacc (E) , amide (F) ...

P3 : hacc (C) , hdon (D) , hacc (E) , ether (F) ...

Structure-based study

Comparative study on the information to include to the background knowledge

	R^2_{cv}
Ligand alone	0.66
Ligand+Protein+ <i>hb/4</i>	0.69
Ligand+<i>water/4+hb/4</i>	0.75
Ligand+Protein+ <i>water/4+hb/4</i>	0.65
FlexX (scoring function)	0.35

- ✓ Need to include receptor information
- ✓ ILP outperforms empirical binding energy predictors
- ✓ Richness of the background knowledge effect on accuracy
- ✓ Closed interpretation of the model as before but
one localized interaction vs. groups involved in few H-bonds

$$\log(1/K_i) = 2.43 + 0.76 \cdot P_1 + 0.91 \cdot P_2 + 0.35 \cdot P_3 - 0.49 \cdot P_4$$

P1 : *hb(C, D), carbonyl(E), amide(F) ...*

P2 : *water(C), alcohol(D), alcohol(E), amide(F) ...*

P3 : *water(C), water(D), alcohol(E), amide(F) ...*

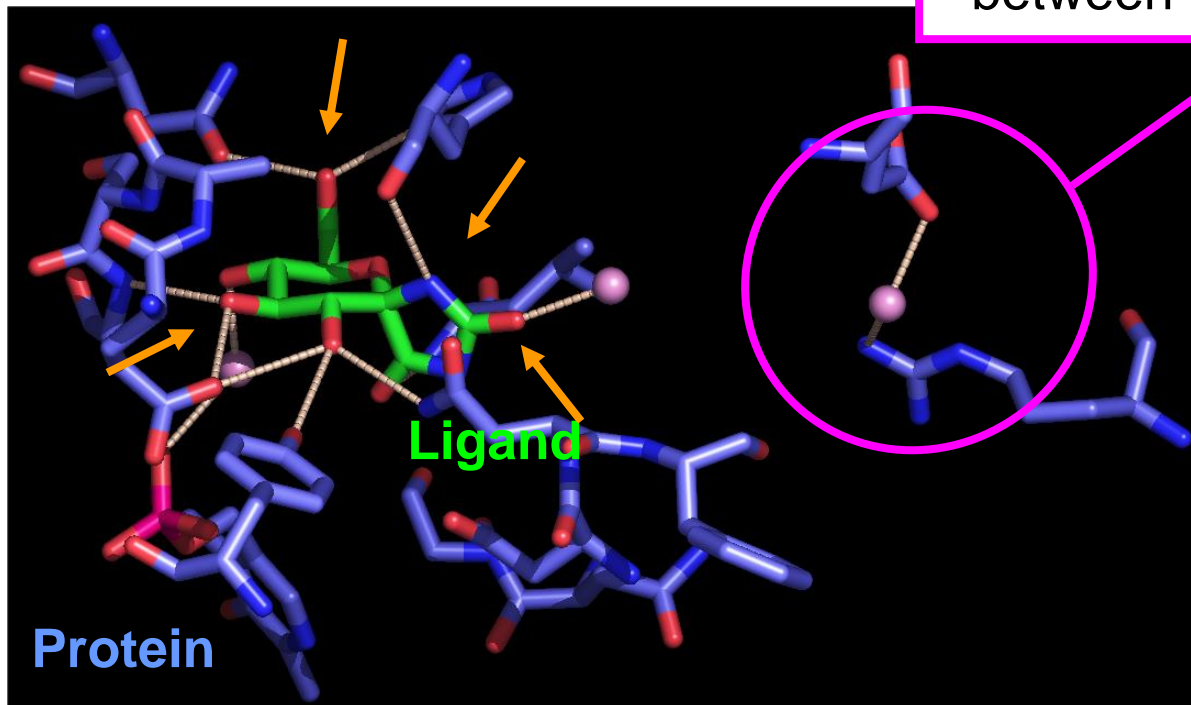
P4 : *water(C), alcohol(D), methylen(E), equiv_ether(F) ...*

Mapping onto the active site

3D representation of the pharmacophores proposed by our Relational model

Groups involved in several non-bonded interactions

Stabilisation of the active site via a water bridge between two amino acids



Enot, D. & King,
R.D. (2003) PKDD
(best paper)

Quantum Mechanics

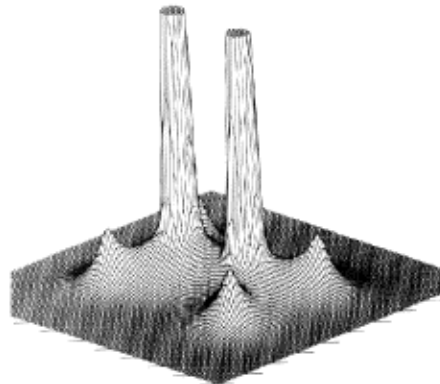
StruQT: Structure representation using quantum topology

Wish-list for molecular structure representations

- n Firmly rooted in quantum mechanics - molecules are quantum mechanical objects not classical ones.
- n Encode efficiently the electronic structure of the molecule.
- n Be related to existing chemical knowledge and concepts.
- n Be efficient for machine inference.

Electron density

- n Electron density is intuitive for the chemist.
- n It is deeply rooted in quantum mechanics.
- n CoMFA methods use 3D grid sampling.



A different approach

- n The 3D density distribution itself is not well suited for machine representation.
- n Use instead the topology of the electron density by inspecting the

$$\nabla\rho(r)$$



Richard F.W. Bader

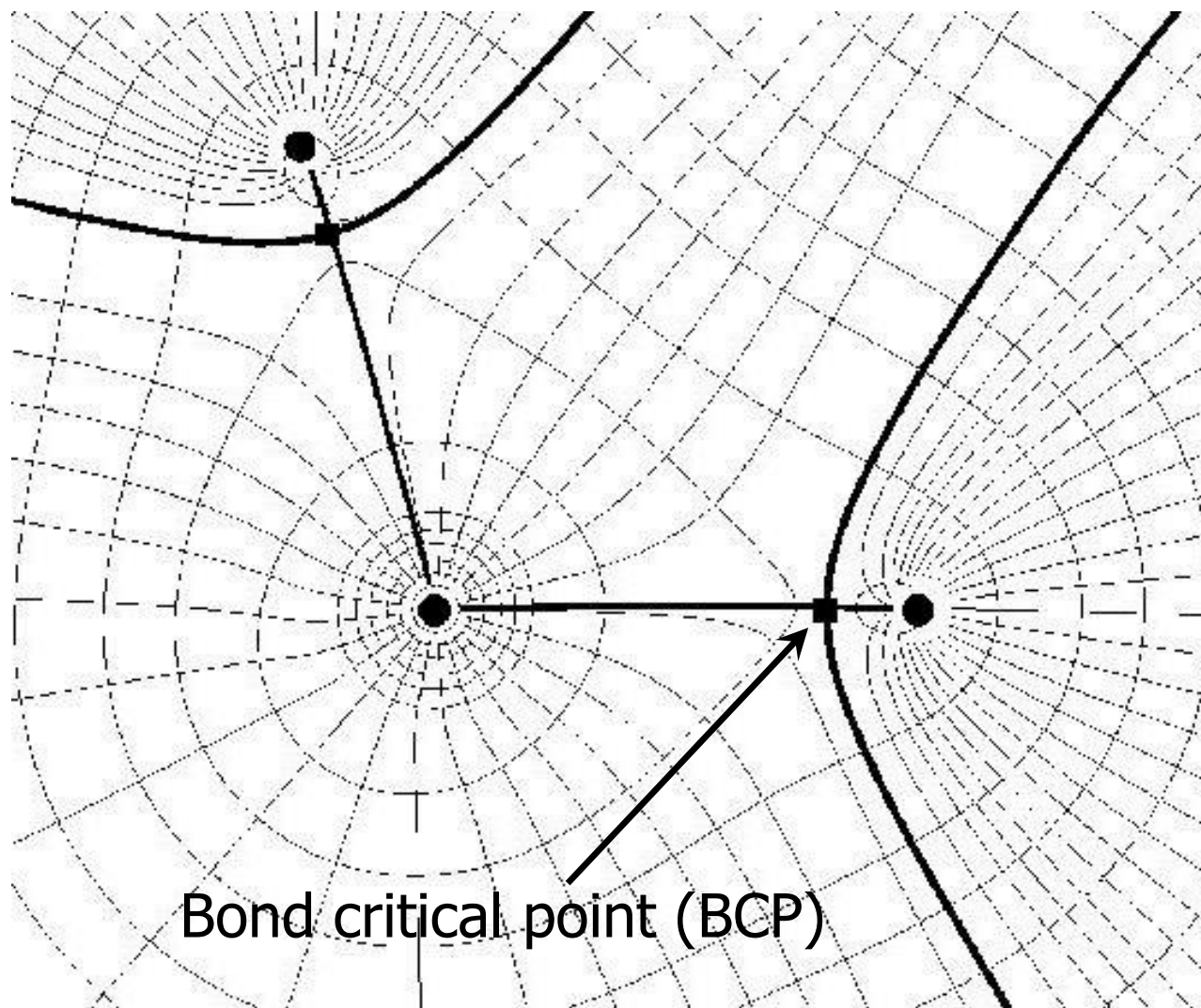
Critical points

- n The topology of the electron density is described by a set of critical points (CP) and their properties
- n CP:

$$\nabla \rho(r) = 0$$

$$\nabla \rho(r) = u_x \frac{\partial \rho}{\partial x} + u_y \frac{\partial \rho}{\partial y} + u_z \frac{\partial \rho}{\partial z}$$

Gradient paths for water



Example StruQT predicates

n	cp(Id,Field,[x,y,z],Hessian)	% critical point
n	signature(Id,Signature).	% signature of Hessian
n	rank(Id,Rank).	% rank of Hessian
n	con_path(Id_1,Id_2)	% connectivity of points
n	hessian(Id,L_1,L_2,L_3)	% eigenvalues of Hessian
n	nuc_attr(Id).	% nuclear attractor
n	bcp(Id)	% bond critical point
n	rcp(Id).	% ring critical point
n	ccp(Id).	% cage critical point

Data Mining

Data Mining

- n The science of finding all frequent “interesting” patterns in data.
- n Not to be confused with machine learning or simply searching databases.
- n The trick is that if you make a pattern more specific then it can only become less frequent. Therefore can deal with vast amounts of data.

Warmr

- n Simplest form: find all frequent combination of items in sets – “basket analysis”. Attribute based.
- n Classical example is finding that beer and nappies are often bought together in supermarkets.
- n Warmr is a Relational algorithm for finding frequent patterns in logical programs.
- n Applied to chemical structures – specific type of program.

Conclusions

- n Finding a good representation is central to learning.
- n QSAR is in essence a learning task.
- n Almost all QSAR methods assume an attribute based representation. This may work well in some circumstances, but it is not a general solution.

- n A more appropriate representation is to use *relations*.
- n The use of relations implies a Relational based learning approach.
- n This has been shown to be effective in many cases.

Acknowledgements

- n David Enot
- n Bjørn Alsberg
- n Nathalie Geneste
- n Luc Dehaspe

BBSRC, EPSRC

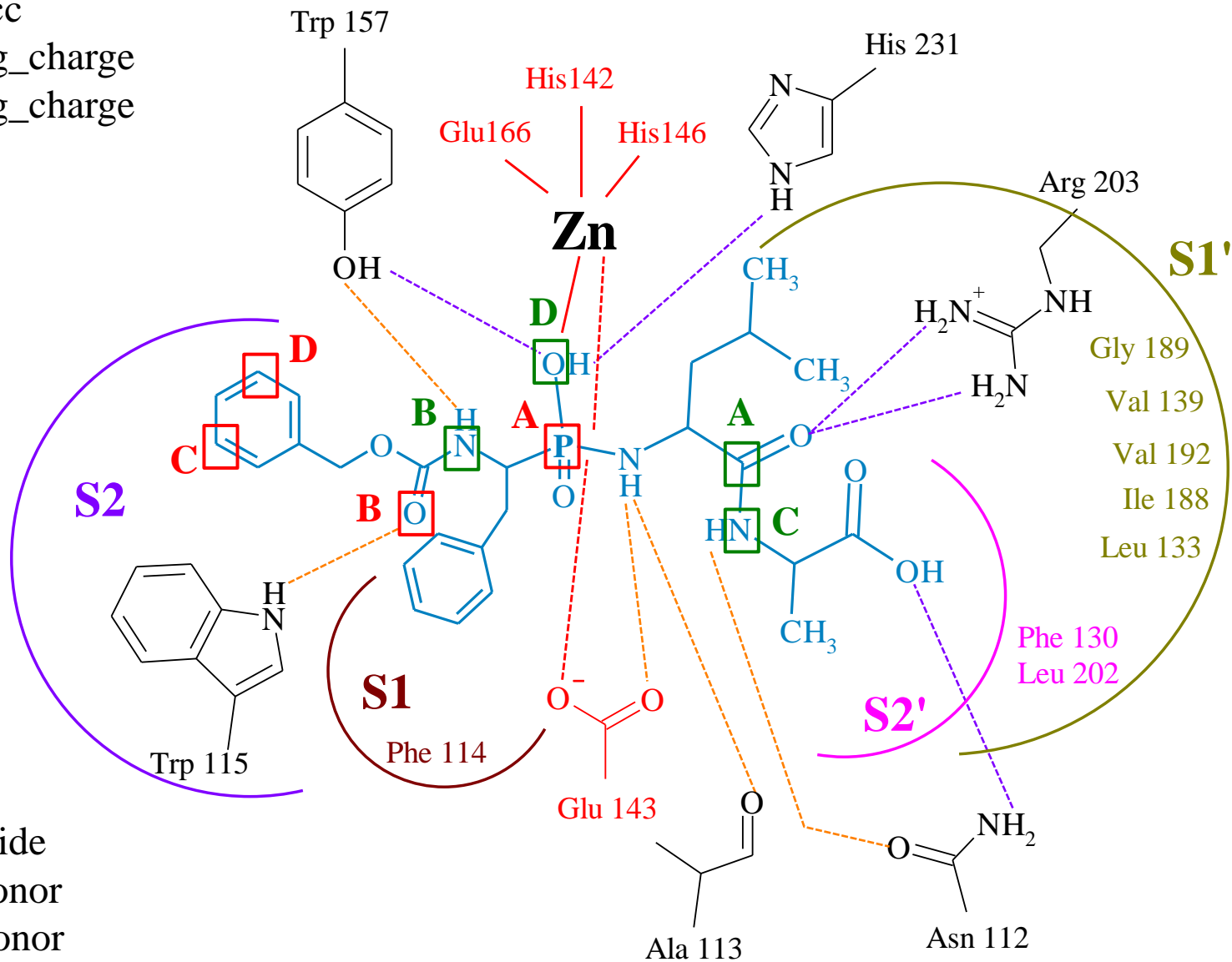
Thermolysin 4TMN

A : phosphorus

B : hacc

C : neg_charge

D : neg_charge



A : amide

B : hdonor

C : hdonor

D : neg_charge

4TMN

