

Handbook of Molecular Descriptors

Roberto Todeschini and Viviana Consonni

 **WILEY-VCH**

Weinheim · New York · Chichester · Brisbane · Singapore · Toronto

This Page Intentionally Left Blank

Handbook of Molecular Descriptors

Roberto Todeschini and Viviana Consonni

 **WILEY-VCH**

Methods and Principles in Medicinal Chemistry

Edited by
R. Mannhold
H. Kubinyi
H. Timmerman

Editorial Board

G. Folkers, H.-D. Höltje, J. Vacca,
H. van de Waterbeemd, T. Wieland

Handbook of Molecular Descriptors

Roberto Todeschini and Viviana Consonni

 **WILEY-VCH**

Weinheim · New York · Chichester · Brisbane · Singapore · Toronto

Series Editors:

Prof. Dr. Raimund Mannhold
Biomedical Research Center
Molecular Drug Research Group
Heinrich-Heine-Universität
Universitätsstraße 1
40225 Düsseldorf
Germany

Prof. Dr. Hugo Kubinyi
Combinatorial Chemistry
and Molecular Modelling
ZHF/G, A 30
BASF AG
67056 Ludwigshafen
Germany

Prof. Dr. Hendrik Timmerman
Faculty of Chemistry
Dept. of Pharmacochimistry
Free University of Amsterdam
De Boelelaan 1083
1081 HV Amsterdam
The Netherlands

Authors:

Prof. Dr. R. Todeschini
Dr. V. Consonni
Milano Chemometrics and QSAR Research Group
Dip. di Scienze dell'Ambiente e del Territorio
Universita degli Studi Di Milano-Bicocca
Piazza della Scienza 1
20126 Milano
Italy

This book was carefully produced. Nevertheless, authors, editors and publisher do not warrant the information contained therein to be free of errors. Readers are advised to keep in mind that statements, data, illustrations, procedural details or other items may inadvertently be inaccurate.

Cover: Staatliche Museen zu Berlin – Preußischer Kulturbesitz
Vorderasiatisches Museum

Library of Congress Card No. applied for

British Library Cataloguing-in-Publication Data: A catalogue record for this book is available from the British Library

Die Deutsche Bibliothek – CIP-Cataloguing-in-Publication Data:
A catalogue record for this book is available from Die Deutsche Bibliothek

ISBN 3-52-29913-0

© WILEY-VCH Verlag GmbH, D-69469 Weinheim (Federal Republic of Germany), 2000

Printed on acid-free paper.

All rights reserved (including those of translation in other languages). No part of this book may be reproduced in any form – by photoprinting, microfilm, or any other means – nor transmitted or translated into machine language without written permission from the publishers. Registered names, trademarks, etc. used in this book, even when not specifically marked as such, are not to be considered unprotected by law.

Composition: Kühn & Weyh, D-79111 Freiburg
Printing: Betzdruck GmbH, D-63291 Darmstadt
Bookbinding: Osswald & Co., D-67433 Neustadt (Weinstraße)

Printed in the Federal Republic of Germany.

*To our loved ones
R.T. and V.C.*

Any alternative viewpoint with a different emphasis leads to an inequivalent description. There is only one reality but there are many viewpoints. It would be very narrow-minded to use only one: we have to learn to be able to imagine several.

Hans Primas

*In: Chemistry, Quantum Mechanics and Reductionism
(Springer-Verlag, 1981)*

Cover

Dragon
Babylon Iŝtar Gate (600 – 500 B.C.)
(Pergamon Museum of Berlin)

A molecular descriptor can be thought of as a mythological animal having several different meanings which depend on one's point of view.

Preface

In the late 1930s, the Hammett equation marked a breakthrough in the understanding of organic chemistry. It describes rate and equilibrium constants of the reactions of aromatic acids, phenols and anilines, as well as other compounds, in a quantitative manner, by using reaction and substituent parameters, ρ and σ . In the same manner, the lipophilicity parameter π , derived from *n*-octanol/water partition coefficients by Corwin Hansch, led to a breakthrough in quantitative structure-activity relationships (QSAR) in biology. Like σ , π also is an additive constitutive molecular property. Their combination, later also with molar refractivity values MR, Taft's steric E_s values or the Verloop steric parameters, allowed the derivation of quantitative models for many biological *in vitro* activity values. Nonlinear lipophilicity models describe *in vivo* biological activities, where substance transport through membranes and distribution within the biological systems play an important role.

Practical problems in the estimation of the lipophilicity of araliphatic and aliphatic compounds led to the *f* hydrophobicity scales of Rekker and Leo/Hansch. However, all such descriptor scales depend on experimental determinations. New molecular descriptors were developed from scratch, starting with the work of Randić, Kier and Hall, i.e. the various molecular connectivity parameters χ . Later the electrotopological state parameters and the Todeschini WHIM parameters were added. Whereas topological descriptors are mathematical constructs that have no unique chemical meaning, they are clearly related to some physicochemical properties and are suited to the description of compound similarities in a quantitative manner. Thus, despite several critical comments in the past, they are now relatively widely used in QSAR studies. Only a meaningless and excessive application in quantitative models, as far as the number of tested and included variables is concerned, still deserves criticism.

This book is a long-awaited monograph on the various properties and molecular descriptors that are of importance in studies of chemical, physicochemical, and biological properties. It is a must for every research worker who is active in this field because it provides an encyclopedic overview of all known descriptors, whether they are physicochemical or topological in their nature. An exhaustive list of references points the way to the original literature.

The series editors wish this book a wide distribution. It is up to the reader to find out which of the properties and descriptors might be most suitable for describing the data. However, the early warnings by Corwin Hansch, John Topliss, and others should not be forgotten: make your model as simple as possible; test and include only a few parameters; try to achieve an understanding of your model; use a test set to check the external predictivity of your model. Molecular descriptors are powerful tools in QSAR studies – but their abuse may lead nowhere. May this book further contribute to their selective and proper use!

August 2000

Raimund Mannhold, Düsseldorf
Hugo Kubinyi, Ludwigshafen
Hendrik Timmerman, Amsterdam

This Page Intentionally Left Blank

Contents

Introduction	XI
User's Guide	XV
Notations and Symbols	XVII
Acronyms	XIX
A – Z	1
Greek Alphabet Entries	511
Numerical Entries	513
Appendix A. Greek Alphabet	515
Appendix B. Symbols of Molecular Descriptors	516
Appendix C. Software	521
References	524

This Page Intentionally Left Blank

Introduction

The effort being made today to organize Knowledge is also a way of participating in the evolution of Knowledge itself. In fact, the significance of attempting such organization can be looked for in its ability not only to give information but also to create know-how: it provides not only a collection of facts – a store of information – but also a contribution to the evolution of Knowledge. The fact is that splitting the organization of Knowledge from its production is completely arbitrary: actually, Knowledge organization is itself one way of doing research.

The true end of an encyclopedic guide is to contribute to the growth of knowledge, but not to knowledge given once and for all, based on some final basic theories, but as a *network of models* in progress. A network primarily consists of knots, i.e. objects, facts, theories, statements, and models, the links between the knots being relationships, comparisons, differences, and analogies: such a network of models is something more than a collection of facts, resulting in a powerful engine for analogical reasoning.

With these purposes in mind, the Authors conceived this *Handbook of Molecular Descriptors* as an encyclopedic guide to molecular descriptors.

First, let us look at the definition of molecular descriptor.

The molecular descriptor is the final result of a logical and mathematical procedure which transforms chemical information encoded within a symbolic representation of a molecule into an useful number or the result of some standardized experiment.

Attention is paid to the term “useful” with its double meaning: it means that the number can give more insight into the *interpretation* of the molecular properties *and/or* is able to take part in a model for the *prediction* of some interesting property of other molecules.

Why must we also accept “or”?

A fundamental task is how to predict and understand experimental facts, i.e. physico-chemical properties, biological activities, and environmental behaviour, from symbolic representations of real objects such as molecules by molecular descriptors.

Because of the huge complexity of this problem, it must be understood that during the development of new descriptors, their interpretation can be weak, provisional, or completely lacking, but their predictive ability or usefulness in application to actual problems can be a strong motive for their use. On the other hand, descriptors with poor predictive power can be usefully retained in models when they are well theoretically founded and interpretable because of their ability to encode structural chemical information.

The incompletely realized comprehension of the chemical information provided by molecular descriptors cannot be systematically ascribed to weakness in the descriptors. Actually, our inability to reduce descriptor meanings to well-established chemical concepts is often because newly emergent concepts need new terms in the language and new hierarchically connected levels for scientific explanation. Thus, what is often considered as scientific failure is sometimes the key to new useful knowledge.

In any case, all the molecular descriptors must contain, to varying extents, chemical information, must satisfy some basic invariance properties and general requirements, and must be derived from well-established procedures which enable molecular descriptors to be calculated for any set of molecules. It is obvious – almost trivial – that a single descriptor or a small number of numbers cannot wholly represent the molecular complexity or model all the physico-chemical responses and biological in-

teractions. As a consequence, although we must get used to living with approximate models (*nothing is perfect!*), we have to keep in mind that “approximate” is not a synonym of “useless”.

The field of molecular descriptors is strongly interdisciplinary and involves a mass of different theories. For the definition of molecular descriptors, a knowledge of algebra, graph theory, information theory, computational chemistry, and theories of organic reactivity and physical chemistry is usually required, although at different levels. For the use of the molecular descriptors, a knowledge of statistics, chemometrics, and the principles of the QSAR/QSPR approaches is necessary in addition to the specific knowledge of the problem. Moreover, programming and sophisticated software and hardware are often inseparable fellow-travellers of the researcher in this field.

The historical development of molecular descriptors reflects some of the distinctive characteristics of the most creative scientists, i.e. their capability of being at the same time engaged and/or detached, rational and/or quirky, serious and/or not so serious. Science is a game and the best players appreciate not only the beauty of a discovery by precise and logical reasoning, but also the taste of making a guess, of proposing eccentric hypotheses, of being doubtful and uncertain when confronted by new and complex problems. Molecular descriptors constitute a field where the most diverse strategies for scientific discovery can be found.

Molecular descriptors will probably play an increasing role in scientific growth. In fact, the availability of large numbers of theoretical descriptors containing diverse sources of chemical information would be useful to better understand relationships between molecular structure and experimental evidence, also taking advantage of more and more powerful methods, computational algorithms and fast computers. However, as before, deductive reasoning and analogy, theoretical statements and hazardous hypotheses, determination and perplexity still remain fundamental tools.

The *Handbook of Molecular Descriptors* tries to meet the great interest that the scientific community is showing in this topic. In fact, as well as the solid interest in the quantitative modelling of biological activity, physico-chemical properties, and environmental behaviour of compounds, an increasing interest has been shown by the scientific community in recent years in the fields of combinatorial chemistry, high-throughput screening, substructural analysis, and similarity searching, for which several approaches which are particularly suitable for informatic treatment have been proposed. Thus, several disciplines such as chemistry, pharmacology, environmental sciences, drug design, toxicology, and quality control for health and safety derive great advantages from these methodologies in their scientific and technological development.

Although experimental measurements would be the most direct way of obtaining safe and high-quality information, the time, costs and hazards involved in the experimentation are relevant limiting factors if we wish to entrust our knowledge growth to fully experiment-based strategies. Moreover, and most importantly, experiments are always suggested by theories, and theoretical models are the way we try to understand reality. Therefore, both general and local models play a fundamental role in scientific growth, leading to a better understanding of the properties of studied phenomena.

The *Handbook of Molecular Descriptors* collects the definitions, formulas and short comments of molecular descriptors known in chemical literature, our intention being to consider all the known molecular descriptors. The definitions of technical terms, around 1800 in all, are organized in alphabetical order. The importance of a molecular descriptor definition is not related to its length. Only a few old descriptors, abandoned

or demonstrated as wrong, have been intentionally left out to avoid confusion. An effort was also made to collect bibliographic information appropriate for the work proposed in this Handbook. We are sorry if any relevant descriptor and/or work has been missed out; this has not been done deliberately and we take full responsibility for any omissions.

Many molecular descriptors have been grouped into classes using a mixed taxonomy based on different points of view, in keeping with the leading idea of the Handbook to promote learning by comparison. Descriptors have been often distinguished by their *physico-chemical meaning* or the specific *mathematical tool* used for their calculation.

Some basic concepts and definitions of statistics, chemometrics, algebra, graph theory, similarity/diversity, which are fundamental tools in the development and application of molecular descriptors, are also presented in the Handbook in some detail. More attention has been paid to information content, multivariate correlation, model complexity, variable selection, and parameters for model quality estimation, as these are the characteristic components of modern QSAR/QSPR modelling.

The Handbook contains nothing about the combinatorial algorithms for the generation and enumeration of chemical graphs, the basic principles of statistics, algorithms for descriptor calculations, or experimental techniques for measuring physico-chemical and biological responses. Moreover, relevant chemometric methods such as Partial Least Squares regression (PLS) and other regression methods, classification methods, cluster analysis, and artificial neural networks, which are also widely applied on molecular descriptors, are simply cited; references are given, but no theoretical aspect is presented. Analogously, computational chemistry methods are only quoted as important tools for calculations, but no claim is made here to their detailed explanation.

Information exchange

The Authors would be grateful to all researchers who would like to send their observations and comments on the Handbook contents, information about new descriptors, and bibliographic references. For e-mail submissions address to: *moldes@disat.unimib.it*.

The Authors are activating a website (<http://disat.unimib.it/chm/>) where recent progress in the molecular descriptor field will be reported, tables of calculated descriptor values for standard data sets collected, and links to other related websites proposed. All information received from interested and collaborative researchers will be useful to develop this Internet tool for sharing ideas and experiences on molecular descriptors and QSAR.

Bibliographic references

The reference list covers a period between 1858 and 2000, and includes 3300 references, about 3000 authors and 250 periodicals. The symbol [R] at the end of a reference denotes a publication with a significant list of references.

Acknowledgements

The idea of producing a *Handbook of Molecular Descriptors* was welcomed by several colleagues whom we warmly thank for their suggestions, revisions, bibliographic information, and moral support. We particularly thank Alexander Balaban,

Subhash Basak, Laura Belvisi, Pierre-Alain Carrupt, Claudio Chiorboli, Johann Gasteiger, Paola Gramatica, Peter Jurs, Lemont Kier, Douglas Klein, Hugo Kubinyi, Silvia Lanteri, Alessandro Maiocchi, Marjana Novic, Demetrio Pitea, Lionello Pogliani, Milan Randic, Gianfranco Tantardini, Bernard Testa, Giuseppe Triacchini, and Jure Zupan.

We are also grateful to Bracco SPA of Milan, EniRicerche of Milan, and the National Institute of Chemistry of Ljubljana for support in the bibliographic search.

Roberto Todeschini and Viviana Consonni

Milano, June, 2000

User's Guide

This handbook consists of definitions of technical terms in alphabetical order, each technical term being an *entry* of the Handbook.

Each topic is organized in a hierarchical fashion. By following cross-references (\rightarrow and typeset in italics) one can easily find all the entries pertaining to a topic even if they are not located together. Starting from the topic name itself, one is referred to more and more specific entries in a top-down manner.

Each entry begins with an entry line.

There are three different kinds of entries: *regular*, *referenced*, and *synonym*.

A **regular entry** has its definition immediately after the entry line. A regular entry is typeset in bold face; it is followed by its (ACRONYM and/or SYMBOL), if any, and by its (: *synonyms*), if any. For example:

Wiener index (*W*) (: *Wiener number*)

A **referenced entry** has its definition in the text of another entry indicated by the symbol \rightarrow and typeset in bold face. For example:

Wiener orthogonal operator \rightarrow **algebraic operators**

A **synonym entry** is followed by the symbol “ : ” and its synonym typeset in italics. To find the definition of a synonym entry, if the synonym is a regular entry, one goes directly to the text under the entry line of the synonym word; otherwise, if the synonym is a referenced entry, one goes to the text of the entry indicated by \rightarrow , typeset in bold face letters. For example:

Wiener number : *Wiener index*

walk number : *molecular walk count* \rightarrow **walk counts**

The text of a regular entry may include the definition of one or more referenced entries highlighted in bold face. When there are many referenced entries collected under one regular entry, called a “mega” entry, they are often organized in hierarchical fashion, denoting them by the symbol \bullet . The sub-entries can be in either alphabetic or logical order. For example, in the mega entry “steric descriptors”, the first sub-entries, each followed by the corresponding text, are:

- \bullet **gravitational indices**
- \bullet **Kier steric descriptor**
- \bullet **Austel branching index**

Finally, a referenced entry within a sub-entry has its definition in the text of the sub-entry denoted by the symbol (\odot ...). For example:

WHIM shape \rightarrow **WHIM descriptors** (\odot global WHIM descriptors)


indicates that the index “WHIM shape” is defined in the sub-entry “global WHIM descriptors” of the main entry “WHIM descriptors”.

In the text of a regular entry one is referred to other relevant terms by words in italics indicated by \rightarrow . In order to reach a complete view of the studied topic, we highly recommend reading also the definitions of these words in conjunction with the original entry. For example:

count descriptors

These are simple molecular descriptors based on counting the defined elements of a compound. The most common chemical count descriptors are → *atom number A*, → *bond number B*, → *cyclomatic number C*, → *H-bond acceptor index* and → *H-bond donor index counts*, → *distance-counting descriptors*, → *path counts*, → *walk counts*. ...

Finally, words in italics not indicated by → in the text of a main entry (or sub-entry) denote relevant terms for the topic which are not further explained or whose definition is reported in a successive part of the same entry.

The symbol  at the end of each entry denotes a list of suggested bibliographic references.

We have made a special effort to keep mathematical notation simple and uniform. A collection of the most often appearing symbols are in the next paragraph Notations and Symbols. Moreover, a list of acronyms helps to decipher and locate the full terminologies given in the book.

Notations and Symbols

The notations and symbols used in the Handbook are listed below. In some cases, notations slightly different from those proposed by the Authors are used to avoid confusion with other descriptors and quantities.

Objects

X	molecular descriptor
M	molecule, compound
\mathcal{M}	experimental measure
\mathcal{P}	experimental property
G	graph, molecular graph
MG	multigraph, molecular multigraph

Sets

V	set of vertices of a graph
E	set of edges of a graph
F	set of fragments of a molecule partition
\mathcal{G}	set of points in a 3D grid
${}^m P_{ij}$	set of atoms of the path of order m from the i th to the j th atoms
${}^m \mathcal{P}$	set of paths of order m

Counts

A	number of atoms of a molecule
B	number of bonds of a molecule
C	number of cycles of a molecule
C^+	number of cycles with overlapping of a molecule
G	number of equivalence classes
h_a	number of hydrogens bonded to the atom a
L	principal quantum numbers
n	number of objects, data, molecules
n_x	number of elements with an x -value
N	generic number of elements
N_X	number of atoms, groups, fragments of X-type
M	number of significant principal components or latent variables
p	number of variables
${}^m P$	number of paths of length m
P	total number of paths of a graph
Z_a	atomic number of the atom a

Matrix operators

C	column sum operator
D	diagonal operator
\mathcal{R}	row sum operator
S	total sum operator
\mathcal{W}	Wiener operator

Indices and characteristic symbols

<i>a</i>	index on the atoms of a molecule
<i>b</i>	index on the bonds of a molecule
<i>g</i>	index on the equivalence classes
<i>i, j, k, l, f, m</i>	generic indices
<i>x, y, z</i>	geometric coordinates
<i>d</i>	data distances
<i>d</i>	topological distances
<i>r</i>	geometric distances
δ	vertex degree
δ^b	bond vertex degree
δ^v	valence vertex degree
ϵ	edge degree
η	atom eccentricity
π	bond order
σ	vertex distance degree
<i>m</i>	order of a descriptor, exponent
<i>w</i>	weights, atom properties
<i>p</i>	probability
<i>q</i>	atomic charge
ℓ_{jm}	PCA loadings of the <i>m</i> th component for the <i>j</i> th variable
λ	eigenvalue
$\lambda_j(\mathbf{M})$	<i>j</i> th eigenvalue from the matrix M
$[\mathbf{M}]_{ij}, m_{ij}$	<i>i-j</i> element of the matrix M
t_{im}	<i>i</i> th score of the <i>m</i> th component from PCA or PLS
<i>t_m</i>	<i>m</i> th vector score
<i>v</i>	generic column vector
<i>p_i</i>	<i>i</i> th grid point of coordinates (<i>x, y, z</i>)
<i>a</i>	vector of atoms in a path
D	dimension (0,1,2,3)
<i>D</i>	diameter
<i>R</i>	radius
I	binary or indicator variable
<i>I</i>	information content

Acronyms

The most well-known acronyms used to define research fields, methods, and molecular descriptors are listed below, in alphabetical order.

AAA	Active Analog Approach
AAC	Augmented Atom Codes
AID	Atomic ID number
ANN	Artificial Neural Networks
ATS	Autocorrelation of a Topological Structure
AWC	Atomic Walk Count
BP-ANN	Back-Propagation Artificial Neural Networks
BIC	Bonding Information Content
BID	Balaban ID number
BLOGP	Bodor LOGP
CADD	Computer-Aided Drug Design
CAMD	Computer-Aided Molecular Design
CAMM	Computer-Aided Molecular Modelling
CASE	Computer-Automated Structure Evaluation
CHEMICALC	Combined Handling of Estimation Methods Intended for Completely Automated LogP Calculation
CIC	Complementary Information Content
CID	Connectivity ID number
CLOGP	Calculated LOGP
CoMFA	Comparative Molecular Field Analysis
CoMMA	Comparative Molecular Moment Analysis
CoMSIA	Comparative Molecular Similarity Indices Analysis
COSV	Common Overlap Steric Volume
CP-ANN	Counter-Propagation Kohonen Artificial Neural Networks
CPK	Corey-Pauling-Koltun volume
CPSA	Charged Partial Surface Areas
CR	Continuum Regression
CSA	Cluster Significance Analysis
DARC	Description, Acquisition, Retrieval Computer system
DD	Drug Design
DFT	Density Functional Theory
DG	Distance Geometry
EAID	Extended Adjacency ID number
EC	Extended Connectivity
ECA	Extended Connectivity Algorithm
ECI	Electronic Charge Index
EEVA	Electronic EigenValue descriptors
EVA	EigenValue descriptors
FEVA	First EigenValue Algorithm
FW	Free-Wilson analysis
GA	Genetic Algorithms
GAI	General a_N -Index
GA-VSS	Genetic Algorithms – Variable Subset Selection
GCSA	Generalized Cluster Significance Analysis
GERM	Genetically Evolved Receptor Models

GFA	Genetic Function Approximation
GIPF	General Interaction Properties Function
GOLPE	Generating Optimal Linear PLS Estimations
G-WHIM	Grid-Weighted Holistic Invariant Molecular descriptors
HASL	Hypothetical Active Site Lattice
HBA	Hydrogen Bond Acceptor
HBD	Hydrogen Bond Donor
HQSAR	Hologram QSAR
HFED	Hydration Free Energy Density
HINT	Hydrophatic INTERactions
HOC	Hierarchically Ordered extended Connectivity
HOMO	Highest Occupied Molecular Orbital
HSA	Hydrated Surface Area
HXID	Hu-Xu ID number
IC	neighbourhood Information Content
ILGS	Iterated Line Graph Sequence
ILS	Intermediate Least Squares regression
ISA	Isotropic Surface Area
IVEC	Iterative Vertex and Edge Centricity algorithm
IVS-PLS	Interactive Variable Selection – Partial Least Squares
K-ANN	Kohonen Artificial Neural Networks
KLOGP	Klopman LOGP
LASRR	Linear Aromatic Substituent Reactivity Relationships
LFER	Linear Free Energy Relationships
LSER	Linear Solvation Energy Relationships
LOMO	Lowest Occupied Molecular Orbital
LOVI	LOcal Vertex Invariant
LSER	Linear Solvation Energy Relationships
LUMO	Lowest Unoccupied Molecular Orbital
MCD	MonteCarlo version of MTD
MCI _s	Molecular Connectivity Indices
MCS	Maximum Common Substructure
MEP	Molecular Electrostatic Potential
MFTA	Molecular Field Topology Analysis
MID	Molecular ID number
MLP	Molecular Lipophilicity Potential
MQSI	Molecular Quantum Similarity Indices
MQSM	Molecular Quantum Similarity Measures
MSA	Molecular Shape Analysis
MSD	Minimal Steric Difference
MSG	Molecular SuperGraph
MTD	Minimal Topological Difference
MTI	Molecular Topological Index
MUSEUM	MUtation and SElection Uncover Models
MWC	Molecular Walk Count
NN	Neural Networks
OASIS	Optimized Approach based on Structural Indices Set
OLS	Ordinary Least Squares regression
PAR	Property-Activity Relationships
PCA	Principal Component Analysis

PCR	Principal Component Regression
PELCO	Pérburbation d'un Environnement Limité Concentrique Ordonné
PID	Prime ID number
PLS	Partial Least Squares regression
PSA	Polar Surface Area
QSAR	Quantitative Structure-Activity Relationships
QShAR	Quantitative Shape-Activity Relationships
QSiAR	Quantitative Similarity-Activity Relationships
QSPR	Quantitative Structure-Property Relationships
QSRC	Quantitative Structure/Response Correlations
QSRR	Quantitative Structure-Reactivity Relationships
RID	Ring ID number
RR	Ridge Regression
RBSM	Receptor Binding Site Model
RSM	Receptor Surface Model
SA	Surface Area
SAR	Structure-Activity Relationships
SASA	Solvent-Accessible Surface Area
SAVOL	Solvent-Accessible VOLume
SBL	Smallest Binary Label
SIBIS	Steric Interactions in BIological Systems
SIC	Structural Information Content
SID	Self-returning ID number
SPP	Submolecular Polarity Parameter
SPR	Structure-Property Relationships
SRC	Structure/Response Correlations
SRW	Self-Returning Walk
SWIM	Spectral Weighted Invariant Molecular descriptors
SWM	Spectral Weighted Molecular signals
SWR	StepWise Regression
TI	Topological Index
TIC	neighbourhood Total Information Content
TLP	Topological Lipophilicity Potential
TLSER	Theoretical Linear Solvation Energy Relationships
TMSA	Total Molecular Surface Area
TOSS-MODE	TOpological SubStructure MOlecular DEsign
VFA	Voronoi Field Analysis
VR	Variable Reduction
VS	Variable Selection
VSS	Variable Subset Selection
WHIM	Weighted Holistic Invariant Molecular descriptors
WID	Weighted ID number
WLN	Wiswesser Line-formula Notation
3D-MoRSE	3D-Molecule Representation of Structures based on Electron dif- fraction descriptors

This Page Intentionally Left Blank

A

A_{x1} , A_{x2} , A_{x3} eigenvalue indices → **eigenvalue-based descriptors**

absolute hardness → **quantum-chemical descriptors** (⊙ hardness indices)

acceptor superdelocalizability : *electrophilic superdelocalizability* → **quantum-chemical descriptors**

ACC transforms → **autocorrelation descriptors**

ACGD index → **charged partial surface area descriptors**

acid dissociation constant → **hydrogen-bonding descriptors**

activation hardness → **quantum-chemical descriptors** (⊙ hardness indices)

acyclic graph : *tree* → **graph**

ADAPT approach

A QSAR approach [Jurs *et al.*, 1979; Jurs *et al.*, 1988], implemented in the homonymous software ADAPT (Automated Data Analysis and Pattern Recognition Toolkit), based on the following steps: a) molecular descriptor generation; b) objective feature selection to discard descriptors which contain redundant or minimal information; c) multiple regression analysis by genetic algorithm or simulated annealing variable selection, or computational → *artificial neural networks*.

ADAPT descriptors fall into three general categories: → *topological indices*, → *geometrical descriptors* (including → *principal moments of inertia*, → *volume descriptors* and → *shadow indices*), and → *electronic descriptors* (including partial atomic charges and the → *dipole moment*); moreover, → *molecular weight*, → *count descriptors*, and a large number of → *substructure descriptors* are also generated. In addition, the → *charged partial surface area descriptors* constitute a fourth class of descriptors derived by combining electronic and geometrical information.

Several molecular properties have been modeled by the ADAPT approach, such as biological activities [Henry *et al.*, 1982; Jurs *et al.*, 1983; Jurs *et al.*, 1985; Walsh and Claxton, 1987; Wessel *et al.*, 1998; Eldred and Jurs, 1999a; Eldred *et al.*, 1999b], boiling points [Smeeks and Jurs, 1990; Stanton *et al.*, 1991; Stanton *et al.*, 1992; Egolf and Jurs, 1993a; Egolf *et al.*, 1994; Wessel and Jurs, 1995a; Wessel and Jurs, 1995b; Goll and Jurs, 1999a], chromatographic indices [Anker *et al.*, 1990; Sutter *et al.*, 1997], aqueous solubilities [Dunnivant *et al.*, 1992; Nelson and Jurs, 1994; Sutter and Jurs, 1996; Mitchell and Jurs, 1998a], critical temperature and pressures [Turner *et al.*, 1998], ion mobility constants [Wessel and Jurs, 1994; Wessel *et al.*, 1996], reaction rate constants [Bakken and Jurs, 1999b; Bakken and Jurs, 1999a].

📖 [Stanton and Jurs, 1992] [Egolf and Jurs, 1992] [Russell *et al.*, 1992] [Egolf and Jurs, 1993b] [Engelhardt and Jurs, 1997] [Mitchell and Jurs, 1997] [Johnson and Jurs, 1999] [Goll and Jurs, 1999b]

additivity model : *Free-Wilson model* → **Free-Wilson analysis**

additive-constitutive models → **group contribution methods**

additive model of inductive effect → **electronic substituent constants** (⊙ inductive electronic constants)

adjacency matrix (A) (: vertex adjacency matrix)

Derived from the \rightarrow molecular graph G , the adjacency matrix \mathbf{A} represents the whole set of connections between adjacent pairs of atoms [Trinajstić, 1992]. The entries a_{ij} of the matrix equal one if vertices v_i and v_j are adjacent (i.e. the atoms i and j are bonded) and zero otherwise. The adjacency matrix is symmetric with dimension $A \times A$, where A is the number of atoms and it is usually derived from an \rightarrow *H-depleted molecular graph*.

The i th row sum of the adjacency matrix is called \rightarrow vertex degree δ_i , defined as:

$$\delta_i = \mathcal{R}_i(\mathbf{A}) = \sum_{j=1}^A a_{ij}$$

where \mathcal{R}_i is the \rightarrow row sum operator; it represents the number of σ -bonds of the i th atom.

The **total adjacency index** A_V is the sum of all the entries of the adjacency matrix of a molecular graph, and is twice the \rightarrow bond number B [Harary, 1969a]:

$$A_V = S(\mathbf{A}) = \sum_{i=1}^A \sum_{j=1}^A a_{ij} = \sum_{i=1}^A \delta_i = 2B$$

where S is the \rightarrow total sum operator. Therefore, the number of entries equal to one in the adjacency matrix is $2B$, while the number of entries equal to zero is $A^2 - 2B$; in particular, for acyclic graphs the total number of entries equal to one is $2(A - 1)$ and the number of entries equal to zero is $A^2 - 2(A - 1)$; for monocyclic graphs they are $2A$ and $A^2 - 2A$, respectively. The total adjacency index is sometimes calculated as the half-sum of the adjacency matrix elements.

Simple \rightarrow topological information indices can be calculated on both the equality and magnitude of adjacency matrix elements. Moreover, \rightarrow walk counts and \rightarrow self-returning walk counts which coincide with the spectral moments of the adjacency matrix are calculated by the increasing powers of the adjacency matrix [McKay, 1977; Jiang *et al.*, 1984; Kiang and Tang, 1986; Hall, 1986; Jiang and Zhang, 1989; Jiang and Zhang, 1990; Markovic and Gutman, 1991; Jiang *et al.*, 1995; Markovic and Stajkovic, 1997; Markovic, 1999].

In order to take into account the heteroatoms in the molecule, the **augmented adjacency matrix** was proposed by Randić [Randić, 1991c; Randić, 1991d; Randić and Dobrowolski, 1998b] replacing the zero diagonal entries of the “normal” adjacency matrix with specific values empirically obtained and characterizing different atoms in the molecule. The row sums of this adjacency matrix are \rightarrow local vertex invariants encoding the connectivity of each atom and its atom type; therefore they can be viewed as augmented vertex degrees. The inverse of the square root of the product of the augmented degrees of the vertices incident with a bond is used as bond weight in calculating the \rightarrow weighted path counts.

Other topological matrices are derived from the adjacency matrix, such as \rightarrow atom connectivity matrices, \rightarrow Laplacian matrix and the powers of the adjacency matrix used to obtain walk counts and the corresponding \rightarrow molecular descriptors.

Moreover, the adjacency matrix can be transformed into a **decimal adjacency vector** \mathbf{a}^{10} of A elements each being a local vertex invariant obtained by the following expression [Schultz and Schultz, 1991]:

$$a_i^{10} = (2 \cdot a_{i1})^{A-1} + (2 \cdot a_{i2})^{A-2} + \dots + (2 \cdot a_{iA})^0$$

where a_{ij} is the j th column element of the i th row of the adjacency matrix \mathbf{A} (zero or one). In this way, the information contained in the adjacency matrix is compressed into an A -dimensional vector. For example, a row of the adjacency matrix [0 1 1 1 0] gives a value of 14, obtained as

$$a_i^{10} = (2 \cdot 0)^{5-1} + (2 \cdot 1)^{5-2} + (2 \cdot 1)^{5-3} + (2 \cdot 1)^{5-4} + (2 \cdot 0)^0 = 14$$

The elements of the decimal adjacency vector are integers which were used for \rightarrow *canonical numbering* of molecular graphs [Randic, 1974].

From the decimal adjacency vector, three different indices were proposed as molecular descriptors:

a) the sum of the elements of the \mathbf{a}^{10} vector, i.e.

$$A1 = \sum_{i=1}^A a_i^{10}$$

b) the sum of the linear combination of vertex degrees δ_i each weighted by the corresponding decimal adjacency vector elements a_i^{10} , i.e.

$$A2 = \sum_{i=1}^A \delta_i \cdot a_i^{10}$$

c) the sum of the elements of the A -dimensional vector \mathbf{d} obtained by multiplying the topological \rightarrow *distance matrix* \mathbf{D} by the decimal adjacency vector, i.e.

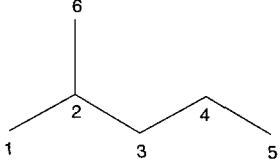
$$A3 = \sum_{i=1}^A [\mathbf{d}]_i$$

where the vector \mathbf{d} is calculated as:

$$\mathbf{d} = \mathbf{D} \cdot \mathbf{a}^{10}$$

Example : 2-methylpentane

adjacency matrix \mathbf{A}



Atom	1	2	3	4	5	6	δ_i
1	0	1	0	0	0	0	1
2	1	0	1	0	0	1	3
3	0	1	0	1	0	0	2
4	0	0	1	0	1	0	2
5	0	0	0	1	0	0	1
6	0	1	0	0	0	0	1

Atom	\mathbf{a}^{10}
1	16
2	41
3	20
4	10
5	4
6	16

$$A_v = \sum_{i=1}^6 \sum_{j=1}^6 a_{ij} = \sum_{i=1}^6 \delta_i = 2 \cdot B = 10$$

$A_1 = 107 \quad A_2 = 219 \quad A_3 = 1028$

Box A-1.

admittance matrix : Laplacian matrix

adsorbability index (AI)

An empirical molecular descriptor proposed by Abe *et al.* (1986) derived from a group contribution method based on molecular refractivity to predict the activated carbon adsorption of 157 compounds from aqueous solutions [Okouchi and Saegusa, 1989]. The adsorbability index for a molecule is calculated by the expression:

$$AI = \sum_i A_i + \sum_i I_i$$

where the sums run over the atoms or functional groups; A indicates the atomic or group factors of increasing or decreasing adsorbability in the molecule and I represents special correction factors accounting for functional group effects.

For example, for benzene: $AI = 6 \cdot A_C + 6 \cdot A_H + 3 \cdot A_{C=C} = 6 \cdot 0.26 + 6 \cdot 0.12 + 3 \cdot 0.19 = 2.85$; for 1,1,2-trichloroethane: $AI = 2 \cdot A_C + 3 \cdot A_H + 3 \cdot A_{Cl} = 2 \cdot 0.26 + 3 \cdot 0.12 + 3 \cdot 0.59 = 2.65$

Table A-1. Values of A and I factors proposed by Abe *et al.*

Atom / Group	A	Group	I
C	0.26	Aliphatic	
H	0.12	–OH (alcohols)	–0.53
N	0.26	–O– (ethers)	–0.36
O	0.17	–CHO (aldehydes)	–0.25
S	0.54	N (amines)	–0.58
Cl	0.59	–COOR (esters)	–0.28
Br	0.86	>C=O (ketones)	–0.30
NO ₂	0.21	–COOH (fatty acids)	–0.03
–C = C–	0.19		
Iso	–0.12	α-Amino acids	–1.55
Tert	–0.32		
Cyclo	–0.28	All groups in aromatics	0

Aihara resonance energy → resonance indices

a_N-index → determinant-based descriptors (⊙ general a_N-index)

algebraic operators

Algebraic operators play a meaningful role in the framework of → *molecular descriptors*, both in deriving molecular descriptors and directly as molecular descriptors.

Let \mathbf{M} be a generic matrix with n rows and p columns, denoted as:

$$\mathbf{M} \equiv [m_{ij}] = \begin{vmatrix} m_{11} & m_{12} & \dots & m_{1p} \\ \vdots & & & \vdots \\ m_{n1} & m_{n2} & \dots & m_{np} \end{vmatrix}$$

The corresponding matrix elements m_{ij} are denoted as:

$$m_{ij} \equiv [\mathbf{M}]_{ij} \equiv (i, j)$$

A column vector \mathbf{v} is a special case of a matrix having n rows and one column; the row vector \mathbf{v}^T is a special case of a matrix having one row and n columns.

Some definitions of matrix algebra [Golub and van Loan, 1983; Mardia *et al.*, 1988], algebraic operators and set theory are given below.

- **characteristic polynomial**

Let \mathbf{M} be a square matrix ($n \times n$) and x a scalar variable, the characteristic polynomial $\mathcal{P}(\mathbf{M}, x)$ is defined as:

$$\mathcal{P}(\mathbf{M}, x) = \det(x\mathbf{I} - \mathbf{M}) = \sum_{k=0}^n a_k x^{n-k}$$

where \mathbf{I} is the identity matrix, i.e. a matrix having the diagonal elements equal to one and all the off-diagonal elements equal to zero, and a_k the polynomial coefficients. Therefore, the characteristic polynomial is obtained by expanding the determinant and then collecting terms with equal powers of x .

The **eigenvalues** λ of the matrix \mathbf{M} are the n roots of its characteristic polynomial and the set of the eigenvalues is called **spectrum of a matrix** $\Lambda(\mathbf{M})$. Determinant and trace of \mathbf{M} are given by the following expressions:

$$\det(\mathbf{M}) = \prod_{k=1}^n \lambda_k \quad \text{tr}(\mathbf{M}) = \sum_{k=1}^n \lambda_k$$

respectively.

For each eigenvalue λ_k , there exists a non-zero vector \mathbf{t}_k satisfying the following relationship:

$$\mathbf{M}\mathbf{t}_k = \lambda_k \mathbf{t}_k$$

The n -dimensional vectors \mathbf{t}_k are called **eigenvectors** of \mathbf{M} .

- **cardinality of a set**

The cardinality of a set S is the number of elements in S and is indicated as $|S|$.

- **column sum operator**

This operator C_j performs the sum of the elements of the j th matrix column:

$$C_j(\mathbf{M}) \equiv \sum_{i=1}^n m_{ij}$$

The **column sum vector**, denoted by \mathbf{cs} , is a p -dimensional vector collecting the results obtained by applying C_j operator on all the p columns of the matrix.

- **determinant**

The determinant of an $n \times n$ square matrix \mathbf{M} , denoted by $\det(\mathbf{M})$, is scalar and is defined as:

$$\det(\mathbf{M}) = \sum_{\pi} s(\pi) \cdot m_{1,i_1} \cdot m_{2,i_2} \cdot \dots \cdot m_{n,i_n}$$

where the summation ranges over all $n!$ permutations π of the symbols $1, 2, \dots, n$. Each permutation π of degree n is given by

$$\pi = \begin{pmatrix} 1 & 2 & \dots & n \\ i_1 & i_2 & \dots & i_n \end{pmatrix}$$

where i_1, i_2, \dots, i_n are the symbols $1, 2, \dots, n$ in some order. The sign function $s(\pi)$ is defined as:

$$s(\pi) = \begin{cases} +1 & \text{if } \pi \text{ is even} \\ -1 & \text{if } \pi \text{ is odd} \end{cases}$$

Related to the definition of determinant are permanent, pfaffian and hafnian.

The **permanent**, denoted by $\text{per}(\mathbf{M})$, also referred to as positive determinant, is defined by omitting the sign function $s(\pi)$ [Schultz *et al.*, 1992; Yang *et al.*, 1994; Cash, 1995a; Cash, 1998] as:

$$\text{per}(\mathbf{M}) = \sum_{\pi} m_{1,i_1} \cdot m_{2,i_2} \cdot \dots \cdot m_{n,i_n}$$

where π runs over the $n!$ permutations.

The **immanant**, denoted by $d_{\lambda}(\mathbf{M})$, is defined as:

$$d_{\lambda}(\mathbf{M}) = \sum_{\pi} \chi_{\lambda}(\pi) \cdot m_{1,i_1} \cdot m_{2,i_2} \cdot \dots \cdot m_{n,i_n}$$

where π runs over the $n!$ permutations. $\chi_{\lambda}(\pi)$ is an irreducible character of the symmetric group indexed by a partition λ of n .

The **pfaffian**, denoted by $\text{pfa}(\mathbf{M})$, is analogous to the determinant where the summation over all the permutations $\pi (i_1, i_2, \dots, i_n)$ must also satisfy the limitations

$$i_1 < i_2, i_3 < i_4, \dots, i_{n-1} < i_n; \quad i_1 < i_3 < i_5 < \dots < i_{n-1}$$

The entries of the main diagonal are excluded from the calculation of the pfaffian [Caianiello, 1953; Caianiello, 1956].

The **hafnian**, denoted by $\text{haf}(\mathbf{M})$, is analogous to the permanent where the summation over all the permutations $\pi (i_1, i_2, \dots, i_n)$ must also satisfy the limitations

$$i_1 < i_2, i_3 < i_4, \dots, i_{n-1} < i_n; \quad i_1 < i_3 < i_5 < \dots < i_{n-1}$$

The entries of the main diagonal are excluded from the calculation of the hafnian.

The hafnian calculated considering only the entries above the main diagonal is called the **short-hafnian**, $\text{shaf}(\mathbf{M})$, while the hafnian calculated considering both entries above and below the main diagonal can also be referred to as the **long-hafnian**, $\text{lhaf}(\mathbf{M})$ [Schultz and Schultz, 1992]. Hafnians and pfaffians differ only in the sign function $s(\pi)$ included in the definition of pfaffian.

For example, for a matrix \mathbf{M} of order 4, pfaffian, long-hafnian and short-hafnian are the following:

$$\text{pfa} = m_{12} \cdot m_{34} - m_{13} \cdot m_{24} + m_{14} \cdot m_{23}$$

$$\text{shaf} = m_{12} \cdot m_{34} + m_{13} \cdot m_{24} + m_{14} \cdot m_{23}$$

$$\text{lhaf} = m_{12} \cdot m_{21} \cdot m_{34} \cdot m_{43} + m_{13} \cdot m_{31} \cdot m_{24} \cdot m_{24} + m_{14} \cdot m_{41} \cdot m_{23} \cdot m_{32}$$

Some molecular descriptors, called \rightarrow *determinant-based descriptors*, are calculated as the determinant of a \rightarrow *matrix representation of a molecular structure*. Moreover, permanents, short- and long-hafnians, calculated on the topological \rightarrow *distance matrix* \mathbf{D} , were used as graph invariants by Schultz and called **per(D) index**, **shaf(D) index**, **lhaf(D) index** [Schultz *et al.*, 1992; Schultz and Schultz, 1992].

📖 [Schultz *et al.*, 1993; Schultz and Schultz, 1993; Schultz *et al.*, 1994; Schultz *et al.*, 1994; Schultz *et al.*, 1995; Schultz *et al.*, 1996; Chan *et al.*, 1997]

• diagonal matrix

A diagonal matrix \mathbf{M} is a square matrix whose diagonal terms m_{ii} are the only nonzero elements. The **diagonal operator** $\mathcal{D}(\mathbf{M})$ is an operator which transforms a generic square matrix \mathbf{M} into a diagonal matrix:

$$\mathcal{D}(\mathbf{M}) = \begin{vmatrix} m_{11} & \dots & 0 & \dots & 0 \\ 0 & \dots & m_{ii} & \dots & 0 \\ 0 & \dots & 0 & \dots & m_{nn} \end{vmatrix}$$