



### (Q)SAR and surroundings...



MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### or alternatively...

Activity Space

		ACUVIL	Activity Space				
		E.P. (A)	E.P. (B)	E. P. (C)	E.P. (D)	E.P. (E)	
<b>Chemical Space</b>	Comp.1	LC <sub>1</sub> (A)					
	Comp.2	<b>LC</b> <sub>2</sub> (A)					
	Comp.3	LC <sub>3</sub> (A)					
	Comp.4	<b>LC</b> <sub>4</sub> (A)					
	Comp.n	LC <sub>n</sub> (A)					

#### **Screening for the specific endpoint**

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

We dim	will ension	discus of this	s late s table	r abo ( <i>n</i> x <i>m</i> )	ut the !!!
	Activit	y Space			
	E.P. (A)	E.P. (B)	E. P. (C)	E.P. (D)	E.P. (E)
Comp.1	LC <sub>1</sub> (A)				
Comp.2	LC <sub>2</sub> (A)				
Comp.3	LC <sub>3</sub> (A)				
Comp.4	<b>LC</b> <sub>4</sub> (A)				
Comp.n	LC <sub>n</sub> (A)				

#### **Screening for the specific endpoint**

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

**Chemical Space** 



## (Q)SAR it is not necessary to remember this definition:



#### Activities (IC<sub>50</sub>, $\mu$ M)

# Any in vivo or in vitro data is affected by both PHARMACODYNAMIC and PHARMACOKINETCS properties of the specific assay.

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

Please: don't start any (quantitative) structureactivity relationship if you're not confident of *quality* of your data activity!!!



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### Please: don't start any (quantitative) structureactivity relationship if you're not confident of *quality* of your data activity!!!



#### Credits: https://en.wikipedia.org/wiki/Accuracy\_and\_precision

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

#### We can start we are real example:

The work reported from The Sandoz Institute for Medical Research on the development of novel analgesic agents (1) can be used as an example of a simple QSAR. In this study, vanillylamides and vanillylthioureas related to capsaicin were prepared and their activity was tested in an in vitro assay which measured  $^{45}Ca^{2+}$  influx into dorsal root ganglia neurons. The data, which was reported as the EC<sub>50</sub>(µM), is shown in Table 1 (note that compound 6f is the most active of the series).

X	Cmpd number	Cmpd name	X	EC <sub>50</sub> (μΜ)
	1	6a	н	11.80 ± 1.90
	2	6b	CI	1.24 ± 0.11
Ö Ö	3	6d	NO <sub>2</sub>	4.58 ± 0.29
	4	6e	CN	26.50 ± 5.87
	5	6f	$C_6H_5$	0.24 ± 0.30
но	6	6g	$N(CH_3)_2$	4.39 ± 0.67
	7	6h	I	0.35 ±0.05
	8	<b>6</b> i	NHCHO	1

1. Christopher S.J. Walpole, Roger Wrigglerworth, Stuart Bevan, Elizabeth A. Campbell, Andy Dray, Iain F. James, Kay J. Mason, martin N. Perkins and Janet Winter, J. Med. Chem., 36, 2381 (1993).

MS	Confidential and Property of ©2012 Molecular Modeling Section
	Dept. Pharmaceutical and Pharmacological Sciences - University of Padova - Ital



#### (Q)SAR: follow me in this wonderful experience

Activity (A<sub>i</sub>)



# Scatterplot... an interesting place where scouting for patterns!!!

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

S.MORO - IA@DSF QSAR \_1

00



### Starting from our table, what is the <u>first</u> <u>descriptor</u> that we can use?



### What is your feeling (

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### Starting from our table, what is the <u>first</u> <u>descriptor</u> that we can use?



# Well, what we can do now... we can calculate the average, of course!

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

## How useful could be the average value to accurately predict the EC50 of cmpd 8?



## The average, 7.24, provides a guess for the value of compound 8 but, how good is this guess?

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

## Do you remember the standard deviation concept? Just a very quick refresh!!!

Credits: https://en.wikipedia.org/wiki/Standard\_deviation



The standard deviation of the data, s, shows how far the activity values are spread about their average. This value provides an indication of the quality of the guess by showing the amount of variability inherent in the data. In our case 9.48!!!

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

## Do you remember the standard deviation concept? Just a very quick refresh!!!

Credits: https://en.wikipedia.org/wiki/Standard\_deviation

Standard Dev		
Population	Sample	
$\sigma = \sqrt{\frac{\sum (X - \mu)^2}{N}}$ $X - \text{The Value in the data distribution}$ $\mu - \text{The population Mean}$ $N - \text{Total Number of Observations}$	$s = \sqrt{\frac{\sum (X - \bar{x})^2}{n - 1}}$ <i>X</i> – The Value in the data distribution $\bar{x}$ – The Sample Mean <i>n</i> - Total Number of Observations	

# Why n - 1? The n-1 (*called Bessel's correction*) in the sample variance formula is a correction factor to adjust the result because the sample mean tends to underestimate the population mean.

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



## We must change the nature of the <u>descriptor</u> but how?



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

### Spencer M. Free and James W. Wilson: A Mathematical Contribution to Structure

Journal of Medicinal Chemistry

Copyright 1964 by the American Chemical Society

Volume 7, Number 4

JULY 6, 1964

#### A Mathematical Contribution to Structure-Activity Studies

Spencer M. Free, Jr., and James W. Wilson

Research and Development Division, Smith Kline and French Laboratories, Philadelphia, Pennsylvania

Received February 4, 1964

A mathematical technique is suggested as a means of describing structure-activity relationships of a series of chemical analogs. The data requirements included specific side chain arrangements and performance characteristics of all analogs tested. Two examples illustrate the use of the additive mathematical model where the performance characteristics are measures of biological activity. The results rank the structural changes per position by estimating the amount of biological response attributed to each change. The estimates are both positive and negative. Several uses for the mathematical solution are suggested.

#### Credits: https://pubs.acs.org/doi/10.1021/jm00334a001

<b>MS</b>	Confidential and Property of ©2012 Molecular Modeling Section	
	Dept. Pharmaceutical and Pharmacological Sciences - University of Padova - Ital	у

Organic chemists who study analog series relate differences in structure to performance characteristics of each compound. The introduction of electronic computers offers opportunities to enhance this effort. Information retrieval techniques provide chemists with selected lists of analogs associated with biological response data. When such data have not been developed through first hand experience, and/or a large number of compounds are available, the determination of structure-activity relationships is more difficult. When the available data meet a limited number of restrictions, mathematical techniques can supplement the organic chemist's intuition.

#### Credits: https://pubs.acs.org/doi/10.1021/jm00334a001

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

The example is presented to illustrate the additive property of the analogs. The mathematical models described in this paper will be based upon the assumption that there is some such "additivity" in series of analogs.

Continuing with the example to illustrate the mathematical model one writes a formula as follows

response = average + effect of R<sub>1</sub> substituent + effect of R<sub>2</sub> substituent



#### Credits: https://pubs.acs.org/doi/10.1021/jm00334a001

MS	Confidential and Property of ©2012 Molecular Modeling Section
	Dept. Pharmaceutical and Pharmacological Sciences - University of Padova - Italy



#### How it Works:



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



*How it Works:*  $\triangle EC_{50} : H \implies R$ 

Cmpd number	Cmpd name	X	EC <sub>50</sub> (μΜ)	ΔEC <sub>50</sub> (μΜ)
1	6a	н	11.80 ± 1.90	0
2	6b	CI	1.24 ± 0.11	- 10.56
3	6d	NO <sub>2</sub>	4.58 ± 0.29	- 7.22
4	6e	CN	26.50 ± 5.87	14.7
5	6f	$C_6H_5$	0.24 ± 0.30	- 11.56
6	6g	N(CH <sub>3</sub> ) <sub>2</sub>	4.39 ± 0.67	- 7.41
7	6h	I	0.35 ±0.05	- 11.45
8	<b>6</b> i	NHCHO	???	???

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



How we can use this  $\triangle EC_{50}$ 

a) Unfortunately, compound 8!!!



estilate the activity of

b) We assume that the biological activity of a molecule can be represented as follow:

### $EC_{50}(R) = EC_{50}(H) + \triangle EC_{50}(R)$

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



How we can use this  $\triangle EC_{50}$ 

a) Following Free-Wilson approach in the case of at least two (or more) substituent, such as in this example:



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



How we can use this  $\triangle EC_{50}$ 

a) In this case, it can assume that the biological activity of a molecule can be represented as follow:

### $EC_{50}(R) = EC_{50}(H) + \triangle EC_{50}(R1) + \triangle EC_{50}(R2)$

b) If it can be assumed that the contribution of the substituent is independent from where is located, it is very easy to anticipate the biological activity of new disubstituited analogs:



or generalizing:

### Activity = $\sum a_i x_i + C$

*a<sub>i</sub>* are the regression coefficients for each substituent;

 $x_i$  are binary variables (1 if the substituent is present, 0 if absent);

C is a constant (could you read the meaning of the costant?)

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### Activity = $\sum a_i x_i + C$

Despite the descriptive simplicity of the additive (linear) model proposed by Free and Wilson, there are rare examples of quantitative structure-activity relationships that have been shown to be describable with acceptable predictive accuracy!





#### Is the world made of straight lines?



### ... at that of the Free-Wilson time it seemed so!!!

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



# I like to introduce Louis Hammett... and his sigma ( $\sigma$ )



Louis Plack Hammett (April 7, 1894 — February 9, 1987)

1940 *Physical Organic Chemistry.* New York; McGraw-Hill.

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### From the birth of organic physical chemistry... to an

96

Louis P. Hammett

Vol. 59

[CONTRIBUTION FROM THE DEPARTMENT OF CHEMISTRY, COLUMBIA UNIVERSITY]

#### The Effect of Structure upon the Reactions of Organic Compounds. Benzene Derivatives

#### By LOUIS P. HAMMETT

The effect of a substituent in the meta or para position of the benzene ring upon the rate or upon the equilibrium of a reaction in which the reacting group is in a side chain attached to the ring may be represented by a simple formula which is valid within a reasonable precision in a surprising variety of cases. The formula is

$$-RT \ln K + RT \ln K^{0} = \Delta F = A/d^{2} \left( \frac{B_{1}}{D} + B_{2} \right)$$

K is a rate constant or an equilibrium constant for a substituted reactant,  $K^0$  is the corresponding quantity for the unsubstituted reactant,  $\Delta F$  is a free energy change or its kinetic analog, d is the distance from the substituent to the reacting group, D is the dielectric constant of the medium in which the reaction occurs, and the quantities A,  $B_1$  and  $B_2$  are constants independent of temperature and solvent. Of these A depends only upon the substituent and its position in the ring relative to the reacting group (with one exception, the two values necessary for the para nitro group), while  $B_1$  and  $B_2$  depend only upon the reaction.

The most important practical feature of equation (1) is the separation of the effect of a substituent into two constants, one of which depends on

 $\sigma$  is a substituent constant, dependent upon the substituent; p is a reaction constant, dependent upon the reaction, the medium and the temperature. Since the only data available consist of values of the  $\sigma \rho$  product, it is necessary to assign an arbitrary value to some one  $\sigma$  or  $\rho$ . The choice of a value of unity for the  $\rho$  constant in the ionization equilibrium of substituted benzoic acids in water solution at 25° was determined by the large amount of accurate data availabe from the recent work of Dippy and co-workers.\* On this basis the difference between the logarithm of the ionization constant of a substituted benzoic acid and the logarithm of the ionization constant of benzoic acid gives the value of the  $\sigma$  constant for that substituent. With the nucleus of  $\sigma$  values thus provided, p values have been derived by least squares methods for other reactions, and from these in turn  $\sigma$  values have been obtained for substituents whose effects upon the ionization constant of benzoic acid are unknown or inaccurately known. After any new  $\sigma$  value was obtained it was used for the calculation of subsequent p values, so that the order of the calculations, which is that of the key numbers in the Table is of a set of the set of t

L.P. Hammett, The Effect of Structure upon the Reactions of Organic Compounds. Benzene Derivatives. J. Am. Chem. Soc. 1937, 59, 96-103.

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### From the electronic effects of the substituent to the definition of the constant $\sigma$ :

pKa for some substituted benzoic acids (H<sub>2</sub>O, 25°C) and corresponding s values of  $\sigma$  the substituents.

Subtituent X	pKa (XC <sub>6</sub> H₄COOH)	$\sigma_{x}$
Н	4.20	0
m-OCH <sub>3</sub>	4.09	0.11
<i>m</i> -F	3.86	0.34
m-NO <sub>2</sub>	3.49	0.71
p-NO <sub>2</sub>	3.42	0.78
p-CH <sub>3</sub>	4.37	-0.17
p-OCH <sub>3</sub>	4.48	-0.28

 $\sigma_{X} = pKa(C_{6}H_{5}COOH) - pKa(XC_{6}H_{4}COOH)$ 

$$\sigma_x = \log \frac{Ka(XC_6H_4COOH)}{Ka(C_6H_5COOH)}$$

	Confidential and Property of ©2012 Molecular Modeling Section
D	Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

Sostituente	σ <sub>m</sub>	$\sigma_{p}$
NO <sub>2</sub>	0.71	0.78
CN	0.61	0.70
CF <sub>3</sub>	0.43	0.54
CH₃COO	0.39	0.31
Br	0.39	0.23
CH₃CO	0.38	0.48
$CO_2C_2H_5$	0.37	0.45
CI	0.37	0.22
СНО	0.36	0.44
CO <sub>2</sub> H	0.35	0.44
I	0.35	0.28
F	0.34	0.06
C≡CH	0.20	0.23
SCH <sub>3</sub>	0.15	≈ 0
ОН	0.13	-0.38
OCH₃	0.11	-0.28
C <sub>6</sub> H₅	0.05	≈ 0
H	0	0
CH <sub>3</sub>	-0.06	-0.17
C <sub>2</sub> H <sub>5</sub>	-0.07	-0.15
CH(CH <sub>3</sub> ) <sub>2</sub>	-0.07	-0.15
C(CH <sub>3</sub> ) <sub>3</sub>	-0.10	-0.20
N(CH <sub>3</sub> ) <sub>2</sub>	-0.15	-0.63
NH <sub>2</sub>	-0.16	-0.57

#### **Substituent Effect Theory**

To summarize:

+R  $\pi$ -donating; -R  $\pi$ -withdrawing +I  $\sigma$ -donating -I  $\sigma$ -withdrawing

Often the types of substituents are

0.34	0.06		_	_
0.20	0.23	+R, +I	+R, -I	-R, -I
0.15	≈ 0	alkyl,	-NH <sub>2</sub> , -OH, -X, -SH	0 0 0
0.13	-0.38	CH <sub>3</sub> -,	-NR <sub>2</sub> , -OR, -SR	
0.11	-0.28	CH <sub>3</sub> CH <sub>2</sub> -,		
0.05	≈ 0	trialky/silvl-		-C≡N -SO <sub>3</sub> H
0	0	than yionyr,	ö ö	-NO <sub>2</sub> -SO <sub>2</sub> R
-0.06	-0.17			
-0.07	-0.15			
-0.07	-0.15			
-0.10	-0.20			
-0.15	-0.63			
-0.16	-0.57			
ial and Prop	erty of ©2012 M	olecular Modeling Section		

Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

IA(0)USF M2h **5.WORU** 

#### Inductive effect (polarization and field effect)







S.MUKU - IA@USF QSAK 1

MSDept. Pharmaceutical and Property of ©2012 Molecular Modeling Section MSDept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



## ... but is sigma ( $\sigma$ ) an invariant property of the substituent?



MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

#### **Electronic effect of the substituent and constant** $\sigma$ :

### But what electronic effect was Hammett discussing as he defined his constant $\sigma$ ?

$$\sigma_x = \log \frac{Ka(XC_6H_4COOH)}{Ka(C_6H_5COOH)}$$



<b>AS</b>	Confidential and Property of ©2012 Molecular Modeling Section
	Dept. Pharmaceutical and Pharmacological Sciences - University of Padova - Italy

#### Application of the constant $\sigma$ : dissociation of homo-benzoic acids



Dept. Pharmaceutical and Pharmacological Sciences - University of Padova - Italy



## The first example of "additivity" in organic chemistry:

For benzoic acids serie is 1



$$pKa_{(R1,R2)}$$
 =  $pKa_{rif}$  -  $ρΣ σ_{Ri}$ 

a nice and convincing example:

COOH  

$$pKa_{(R1,R2)} = 4.2 - (0.71 \times 2) = 2.78$$
  
 $pKa_{(R1,R2)} = 2.82 (exp)$ 

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy


#### ... and after Hammett

MS	Confidential and Property of ©2012 Molecular Modeling Section
	Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### The fragments method of Hansch-Fujita.

(J. Am. Chem. Soc. 1964, 86, 5175)

Dec. 5, 1964

SUBSTITUENT CONSTANT,  $\pi$ , FROM PARTITION COEFFICIENTS



[CONTRIBUTION FROM THE DEPARTMENT OF CHEMISTRY, POMONA COLLEGE, CLAREMONT, CALIF.]

A New Substituent Constant,  $\pi$ , Derived from Partition Coefficients

By Toshio Fujita, <sup>1a</sup> Junkichi Iwasa, <sup>1b</sup> and Corwin Hansch Received February 19, 1964

The partition coefficients between 1-octanol and water have been determined for 203 mono- and disubstituted benzenes. From these values a substituent constant,  $\pi$ , has been calculated for 67 functional groups. The constant  $\pi$  is defined as:  $\pi = \log P_X - \log P_H$  where  $P_X$  is the partition coefficient of a derivative and  $P_H$  is that of the parent compound.  $\pi$  has been derived for many of the functions from eight different systems: benzene, nitrobenzene, aniline, phenol, benzyl alcohol, benzoic acid, phenylacetic acid, and phenoxyacetic acid. It is found that, although  $\pi$  varies continuously for a given function depending on its electronic environment, the range over which it varies is not great. In certain of the systems,  $\pi$ -values are related by a simple linear expression.



5175



#### **POMONA COLLEGE (Claremont, California)**

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

#### Hydrophobicity and partition coefficient

$$\mu_{(phase2)} = \mu_{(phase2)}^{0} - RT \ln[C_{(phase2)}]$$
Phase 2 (*n*-octanol)
$$\mu_{(phase1)} = \mu_{(phase1)}^{0} - RT \ln[C_{(phase1)}]$$

#### At the equilibrium:

$$\mu_{(phase1)} = \mu_{(phase2)}$$

$$\mu_{(phase1)}^{0} - RT \ln \left[C_{(phase1)}\right] = \mu_{(phase2)}^{0} - RT \ln \left[C_{(phase2)}\right]$$

$$\mu_{(phase1)}^{0} - \mu_{(phase2)}^{0} = RT \ln \left[C_{(phase1)}\right] - RT \ln \left[C_{(phase2)}\right]$$

$$\mu_{(phase1)}^{0} - \mu_{(phase2)}^{0} = RT \ln \left[\frac{C_{(phase1)}}{C_{(phase2)}}\right]$$

$$\left[\frac{C_{(phase1)}}{C_{(phase2)}}\right] = P \text{ Partition coefficient}$$

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

#### Hydrophobicity and partition coefficient

$$\left[\frac{C_{(phasel)}}{C_{(H_2O)}}\right] = P$$
 Partition coefficient

We can define:

*"hydrophobic*" a compound with P > 1; *"hydrophilic*" a compound with P < 1.

$$\log\left[\frac{C_{(phase1)}}{C_{(H_2O)}}\right] = \log P$$

We can define:

#### *"hydrophobic*" a compound with logP > 0; *"hydrophilic*" a compound with logP < 0.



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

#### Hydrophobicity and partition coefficient

#### How we can choose the second phase:



HO

*n*-octanol

•immiscible in water even if 27% of water dissolves in it... so the first region of hydration of the solute is preserved;

- •UV transparent;
- •Low vapor pressure.



cyclohexane

Immiscible in water and very low amount of water dissolves in it... so also the first region of hydration of the solute is lost. The differences between the logP values in *n*octanol and cyclohexane is a measure of the de-hydration energy of a solute.



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### Do you see any similarity? ...in addition to this?

#### How we can choose the second phase:





*n*-octanol

•"Similarity" with biological membrane;



Phospholipid





Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy







Confidential and Property of ©2012 Molecular Modeling Section MSDept. Pharmaceutical and Property of ©2012 Molecular Modeling Section MSDept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy







Cl	0.71
Br	0.86
Ι	1.12
ОН	-0.67
OCH <sub>3</sub>	-0.02
SCH <sub>3</sub>	0.61
CN	-0.57
СООН	-0.28
COOCH <sub>3</sub>	-0.01
COCH <sub>3</sub>	-0.55
NH <sub>2</sub>	-1.23
N(CH <sub>3</sub> ) <sub>2</sub>	-0.28
NO <sub>2</sub>	-0.28
CH <sub>3</sub>	0.56

MSDept. Pharmaceutical and Property of ©2012 Molecular Modeling Section MSDept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy





**Ref. benzene** 

log *P*=2.13

Substituent	$\pi$ aromatic
F	0.14
Cl	0.71
Br	0.86
Ι	1.12
ОН	-0.67
OCH <sub>3</sub>	-0.02
SCH <sub>3</sub>	0.61
CN	-0.57
СООН	-0.28
COOCH <sub>3</sub>	-0.01
COCH <sub>3</sub>	-0.55
$\mathbf{NH}_2$	-1.23
N(CH <sub>3</sub> ) <sub>2</sub>	-0.28
NO <sub>2</sub>	-0.28
CH <sub>3</sub>	0.56



#### Ref. cyclohexane log *P* = 3.44

MSDept. Pharmaceutical and Property of ©2012 Molecular Modeling Section University of Padova - Italy

#### The additive rule:



#### It works? Just check together:







 log P 1.80
 log P 2.30
 log P 2.40

 Clog P 1.61
 Clog P 2.14
 Clog P 2.38

#### ...not too bad!!!

Confidential and Property of ©2012 Molecular Modeling Section MSDept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



But not so good...



MSDept. Pharmaceutical and Property of ©2012 Molecular Modeling Section University of Padova - Italy









#### Corwin... we have a problem!



#### Is this true?

#### But how we can demonstrate this is not true?

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### Corwin... we have a problem!

# Here is how he fixed the problem... as an engineer usually done introduction the magic *correction factors* !!!

Remember : a *correction factor* is any mathematical adjustment made to a calculation to account for deviations in either the sample or the method of measurement.

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### The hydrophobic correction factors, $\Delta \pi$

- Branched carbon chain ( $\Delta \pi$  = -0.20);
- Double bond ( $\Delta \pi$  = -0.30);
- Intra-molecular H-bond ( $\Delta \pi$  = 0.65);
- Ring condensation ( $\Delta \pi = -0.20$ )

#### Finally the log *P* calculated by Hansch-Fujita:

#### Clog P = log P<sub>ref</sub> + $\Sigma \pi_{Xi}$ + $\Sigma \Delta \pi$



#### Nowadays, are these methods reliable?

#### ... you decide!!!



n = 12'202, r<sup>2</sup> = 0.944, r = 0.972, stdev = 0.428, mae = 0.328



#### credits: https://www.molinspiration.com/cgi-bin/properties

#### ClogP... a wonderfully precious molecular descriptor!

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



## Can we replace in our scatterplot the "cmpd number" with " $\sigma_p$ "?

Cmpd number	Cmpd name	X	EC <sub>50</sub> (μΜ)	$\sigma_{p}$
1	6a	н	11.80 ± 1.90	0
2	6b	CI	1.24 ± 0.11	0.78
3	6d	NO <sub>2</sub>	4.58 ± 0.29	0.22
4	6e	CN	26.50 ± 5.87	0.70
5	6f	$C_6H_5$	0.24 ± 0.30	0.10
6	6g	$N(CH_3)_2$	4.39 ± 0.67	-0.63
7	6h	I	0.35 ±0.05	0.28
8	6i	NHCHO	???	0.12

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### Can we replace in our scatterplot the "cmpd number" with " $\sigma_p$ "?



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



## ... or can we replace in our scatterplot the " $\sigma_p$ " with " $\pi$ "?

Cmpd number	Cmpd name	X	EC <sub>50</sub> (μΜ)	σ <sub>p</sub>	π
1	6a	Н	11.80 ± 1.90	0	0
2	6b	CI	1.24 ± 0.11	0.78	0.71
3	6d	NO <sub>2</sub>	4.58 ± 0.29	0.22	- 0.28
4	6e	CN	26.50 ± 5.87	0.70	- 0.57
5	6f	$C_6H_5$	0.24 ± 0.30	0.10	2.00
6	6g	N(CH <sub>3</sub> ) <sub>2</sub>	4.39 ± 0.67	-0.63	- 0.28
7	6h	I.	0.35 ±0.05	0.28	1.12
8	6i	NHCHO	???	0.12	- 0.70

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



## ... or can we replace in our scatterplot the " $\sigma_p$ " with " $\pi$ "?



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### ... and now?

MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### (Q)SAR: we are ready to this...



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



(Q)SAR: we are ready to this...



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### (Q)SAR: we are ready to this...

## DRAGON 7.0 is able to calculate 5270 molecular descriptors.



Padova - Italy S.MOR



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### (Q)SAR: follow me in this wonderful experience



#### Property (P<sub>i</sub>) Scatterplot... an interesting place where scouting for patterns!!!

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### (Q)SAR: how can we select the good "molecular descriptor(s)"?



MSConfidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences - University of Padova - Italy





	Confidential and Property of ©2012 Molecular Modeling Section
<b>C</b> ]	Dept. Pharmaceutical and Pharmacological Sciences - University of Padova - Ital



### Yes, we can look for "*regularity*" (pattern) between the variability of molecular descriptors and the corresponding variability of experimental activities.

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### The beauty of mathematics:



Property (P<sub>i</sub>)

### Discrete (few *x-y* corrispondences) Continuum (~ *x-y* corrispondences)

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



Patterns are gorgeous:

•Patterns can be mathematically condensed in equations; Pattern can be used to describe relationships among variables; Patterns can be used to predict new data; Patterns can be used to verify exiting data;

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### sometimes too gorgeous...



MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### sometimes too gorgeous...

- classification models (qualitative response)
- regression models (quantitative response)
- ranking models (ordered response)

	Confidential and Property of ©2012 Molecular Modeling Section
<b>ID</b>	Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### The scatter plot: the best place where explore (Q)SAR.



Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



## The scatter plot: the best place where explore (Q)SAR.



Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy


# sometimes too gorgeous... the first powerfull example of AI:

- classification models (qualitative response)
- regression models (quantitative response)
- ranking models (ordered response)

	Confidential and Property of ©2012 Molecular Modeling Section
1D	Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

# The scatter plot: the best place where explore (Q)SAR: look at this experiment



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



The two statistical gold rules do build up linear models:

For each independent variable (*molecular descriptor*) you need <u>at</u> <u>least</u> five (5) dependent variable values (*activities*).

The dependent variable values
(*activities*) must be <u>accurate</u> and
<u>precise</u>.

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



MC	Confidential and Property of ©2012 Molecular Modeling Section	
MO	Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Ita	ly



# Now, how can we select the "good" linear model:

### Homoscedasticity

or *homogeneity of variances*, is an assumption of equal or similar variances in different groups (experiments) being compared.



#### credits: https://en.wikipedia.org/wiki/Homoscedasticity\_and\_heteroscedasticity homogeneity of variance.

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

Activity (A<sub>i</sub>)





For homoskedasticity violation, the standard errors will be biased and estimates of regression coefficients will be less efficient. In practice, this usually mean overestimating the precision of your model

credits: https://en.wikipedia.org/wiki/Homoscedasticity\_and\_heteroscedasticity homogeneity of variance.

# The scatter plot: the best place where explore (Q)SAR: look at this experiment



Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



do you remember... Least Squares Analysis?

Credits: https://en.wikipedia.org/wiki/Least\_squares

#### 1. What is a Least Squares Regression Line?

If your data shows a *linear relationship* between the X and Y variables, you will want to find the line that best fits that relationship. That line is called a *Regression Line* and has the equation  $\hat{y}$ = a + b x. The Least Squares Regression Line is the line that makes the vertical distance from the data points to the regression line as small as possible. It's called a "least squares" because the best line of fit is one that minimizes the variance (the sum of squares of the errors).



Activity (A<sub>i</sub>)

### do you remember... Least Squares Analysis?

Credits: https://en.wikipedia.org/wiki/Least\_squares



LSA is a method for linear regression that determines the values of unknown quantities in a statistical model by minimizing the sum of the residuals, the difference between the predicted  $(\hat{y})$  and observed values (y) squared.

**e = ŷ** - y

#### **Property (P<sub>i</sub>)**

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

### do you remember... Least Squares Analysis?

This can be a bit hard to visualize but the main point is you are aiming to find the equation that fits the points as closely as possible.



#### Please looks this video: https://youtu.be/bhKXKTaQ96M

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



**Goodness of fit:** variation in the data is quantified by the *coefficient of determination*  $(r^2)$  which measures how closely the observed data tracks the fitted regression line. Errors in either the model or in the data will lead to a bad fit. This indicator of fit to the regression line is calculated as:



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

### Calculating $r^2$

- Original variance:
- Explained variance:
- Variance around regression line:

$$TSS = \sum_{i=1}^{N} (y_i - \overline{y})^2$$
$$ESS = \sum_{i=1}^{N} (y_{i,calc} - \overline{y})^2$$
$$RSS = \sum_{i=1}^{N} (y_i - y_{calc,i})^2$$

$$r^{2} = \frac{ESS}{TSS} \equiv \frac{TSS - RSS}{TSS} \equiv 1 - \frac{RSS}{TSS} \qquad \qquad \mathbf{0} < \mathbf{r}^{2} < \mathbf{1}$$

Possible values reported for  $r^2$  fall between 0 and 1. For example: with  $r^2$  of 0.83, you can say that 83% of the variability in activity can be explained by the different value of the selected molecular property. *The remaining 17% of variability is due to other unexplained factors.* 

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy **Goodness of fit:** the *Pearson correlation coefficient* (*r*) is the square root of  $r^2$  expressed as a decimal. Its *size* is always between 0 and 1. The *sign* of the correlation coefficient depends on the slope of the regression line:





Property (P<sub>i</sub>) inverse corr.

**Property (P<sub>i</sub>)** 

A perfect correlation of  $\pm$  1 occurs only when the data points all lie exactly on a straight line. A correlation greater than 0.8 would be described as strong, whereas a correlation less than 0.5 would be described as weak.

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



# How we can select the *good* descriptors? *Compare r!*



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

## **Outliers:** an outlier is an observation that is numerically distant from the rest of the data.



### How do we deal with them, usually?



#### **Property (P<sub>i</sub>)**

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### **Be carefull...**

Activity (A<sub>i</sub>)



**Property** (P<sub>i</sub>)

## ... the rabbit (pattern) is out there!!!

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



Acrivity (A<sub>i</sub>)

# The scatter plot: the best place where explore (Q)SAR.



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



# Cross-validation (CV) for detecting and preventing overfitting!

Activity (A<sub>i</sub>)



Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### The "test set" method...looks this:





### The "test set" method...looks this:



1. Randomly choose 30% of the data to be in a *test set*;

2. The remainder is a *training set*;

3. Perform your regression on the *training set*;

4. Estimate your future performance with the *test set*.

#### Mean Squared Error (MSE)

 $(x_i - \overline{x})^2$ n

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



#### Good news:

•Very very simple;

•Can then simply choose the method with the best "test set" score.

#### Bad news:

•Wastes data: we get an estimate of the best method to apply to 30% less data; •If we don't have much data, our test-set might just be lucky or unlucky.

S Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



# or *"LOOCV"* (Leave-One-Out Cross Validation) method... looks this:

Activity (A;)

Property (P<sub>i</sub>)

For each data consider this loop:

- 1. Select the first  $(x_i, y_i)$  data;
- 2. Temporary remove  $(x_i, y_i)$  from the data set;
- 3. Train on the remaining n-1 datapoints;
  - 4. Note your error (x<sub>i</sub>, y<sub>i</sub>);

## When you've done all points, report the mean squared errors (MSE).

Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



MS Confidential and Property of ©2012 Molecular Modeling Section Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



## So... which kind of validation?

	Downside	Upside
Test-set	Variance: unreliable estimate of future performance	Time cheap
Leave- one-out	Time expensive. Has some weird behaviour	Doesn't waste data

	Confidential and Property of ©2012 Molecular Modeling Section
<b>D</b>	Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy



### Here a possible work-flow:

