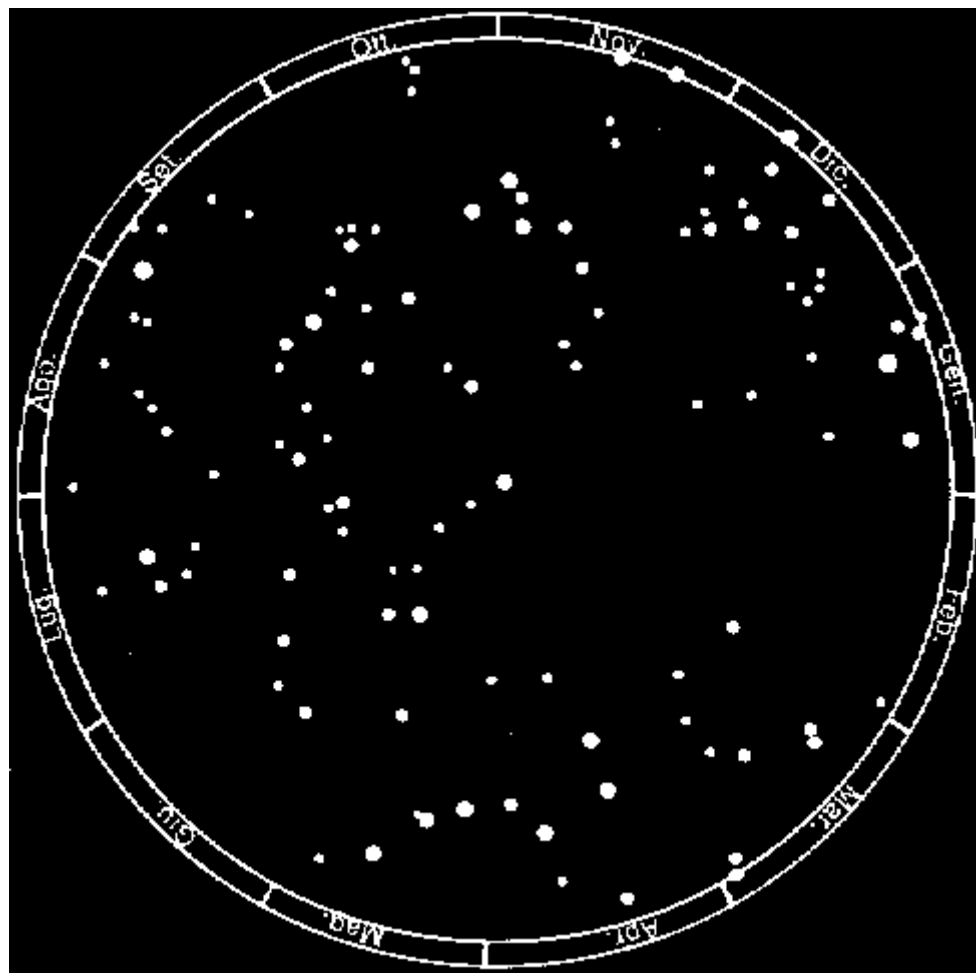
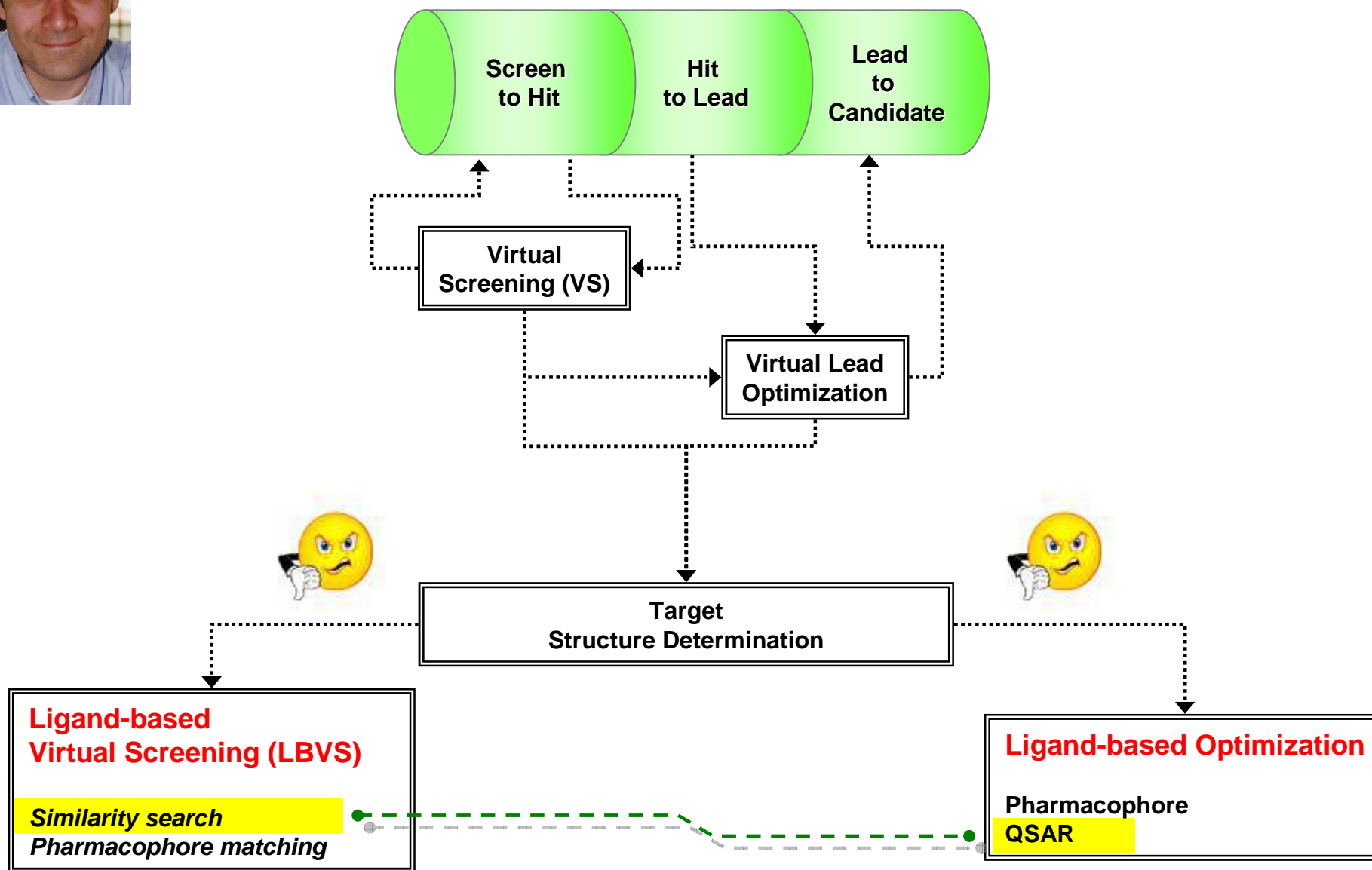


In and around... (Q)SAR





here we are again:





... do you surely remember:

Biological Space

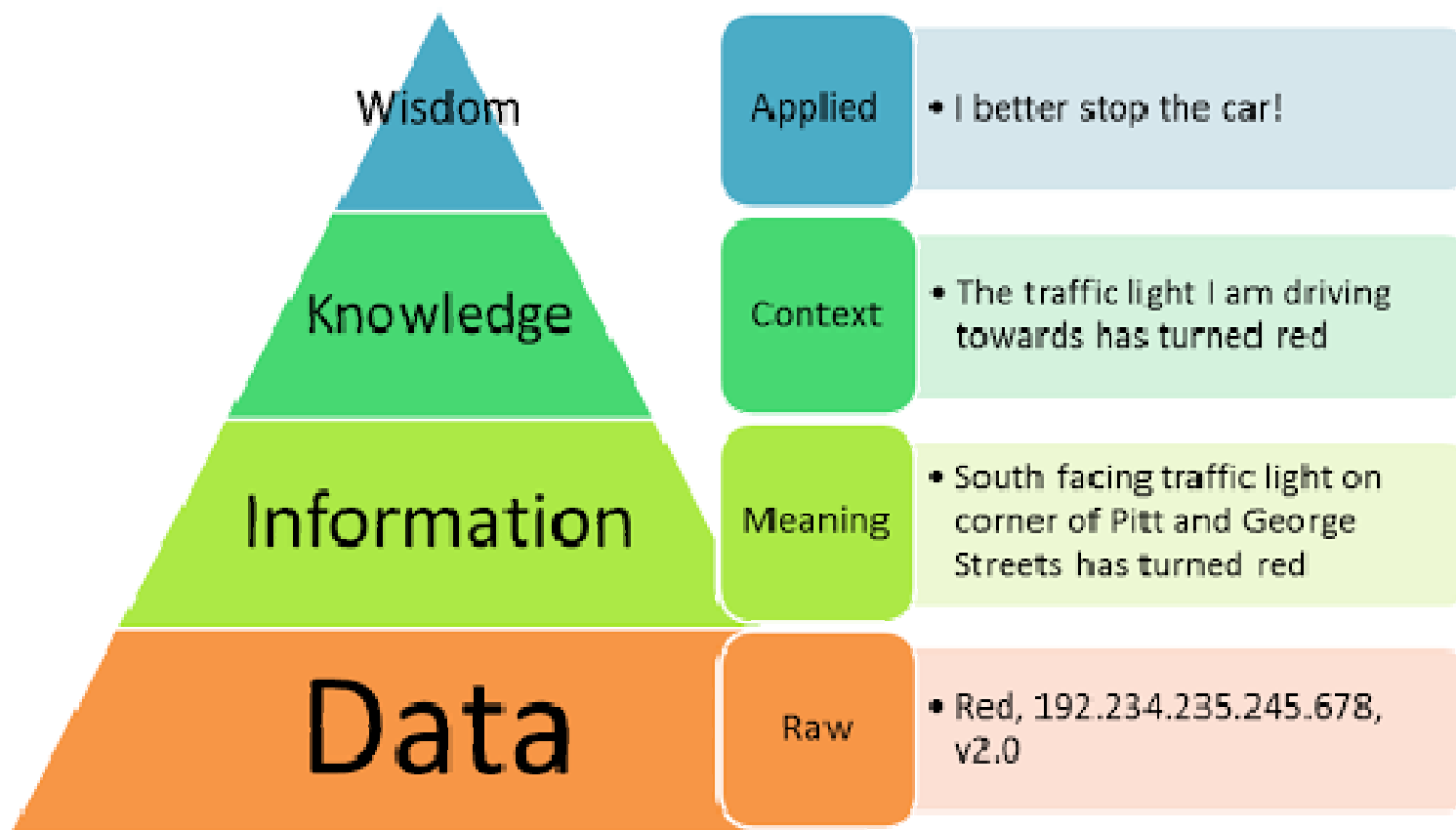
Chemical Space

	TargetA	TargetB	TargetC	TargetD	TargetN
Comp.1	$K_1(A)$				
Comp.2	$K_2(A)$				
Comp.3	$K_3(A)$				
Comp.4	$K_4(A)$				
Comp.n	$K_n(A)$				

Screening



... and surely also this concept:



© 2011 Angus McDonald



Here is another wonderful example of data
→ information transformation:

Biological Space

Chemical Space

	TargetA (EC ₅₀ , μ M)	TargetB	TargetC	TargetD	TargetN
Comp.4	0.5 \pm 0.1				
Comp.3	1.3 \pm 0.1	SORT!			
Comp.1	2.0 \pm 0.1				
Comp.5	3.3 \pm 0.1				
Comp.2	8.3 \pm 0.1				

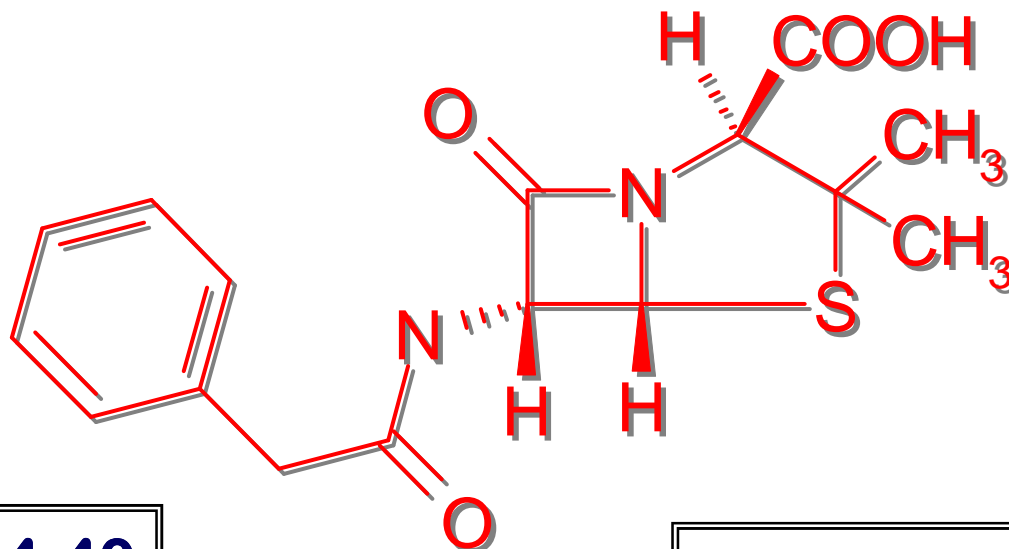


... please: do not start any quantitative structure-activity relationship if you're not confident of 'reproducibility' of your data of biological activity!!!

$$\begin{aligned} EC_{50} &= 1.0 \mu M \neq \\ &1.0 \pm 0.1 \mu M \neq \\ &1.0 \pm 0.5 \mu M \neq \\ &1.0 \pm 1.0 \mu M \end{aligned}$$



NUMB3RS!!!



MW = 334,40

HB_Acc = 4

pKa = 2.7

Volume = 302,37

nC = 16

logP = 0.8

...



We can start filling this table in an other way:

Chemical Space		TargetA (EC ₅₀ , μ M)	Descriptor A	Descriptor B	Descriptor C	Descriptor n
	Comp.4	0.5 \pm 0.1				
	Comp.3	1.1 \pm 0.1				
	Comp.1	2.0 \pm 0.1				
	Comp.5	3.3 \pm 0.1				
	Comp.2	8.3 \pm 0.1				



... how many molecular descriptor can be calculated?

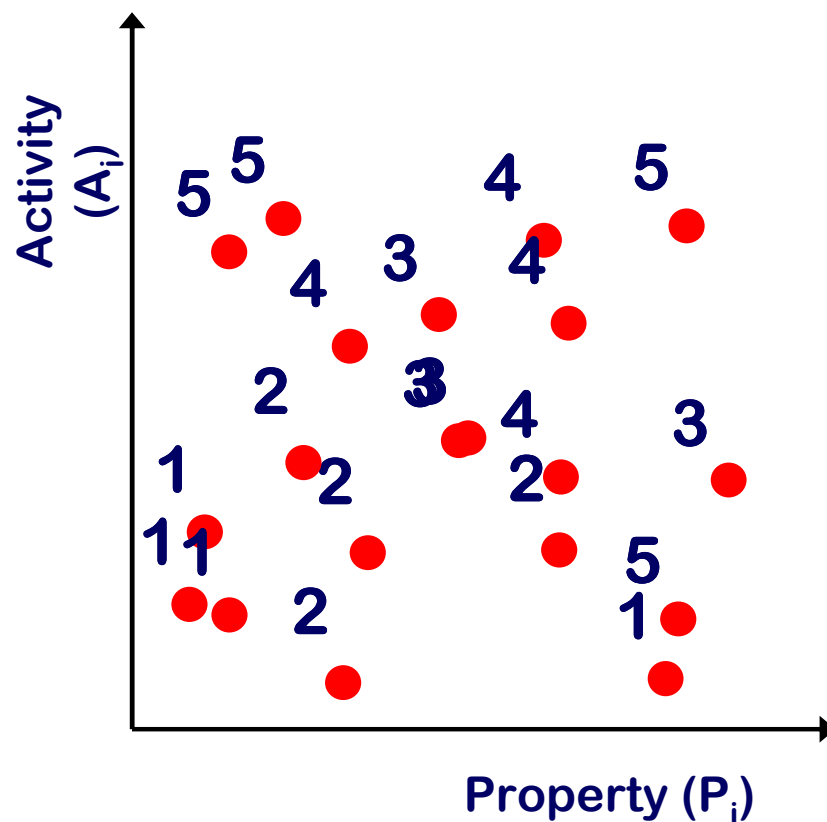
DRAGON 6.0 is able to calculate **4885** molecular descriptors!



<http://www.taletе.mi.it/index.htm>



and so....



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



How can we select the good “*molecular descriptor(s)*”?

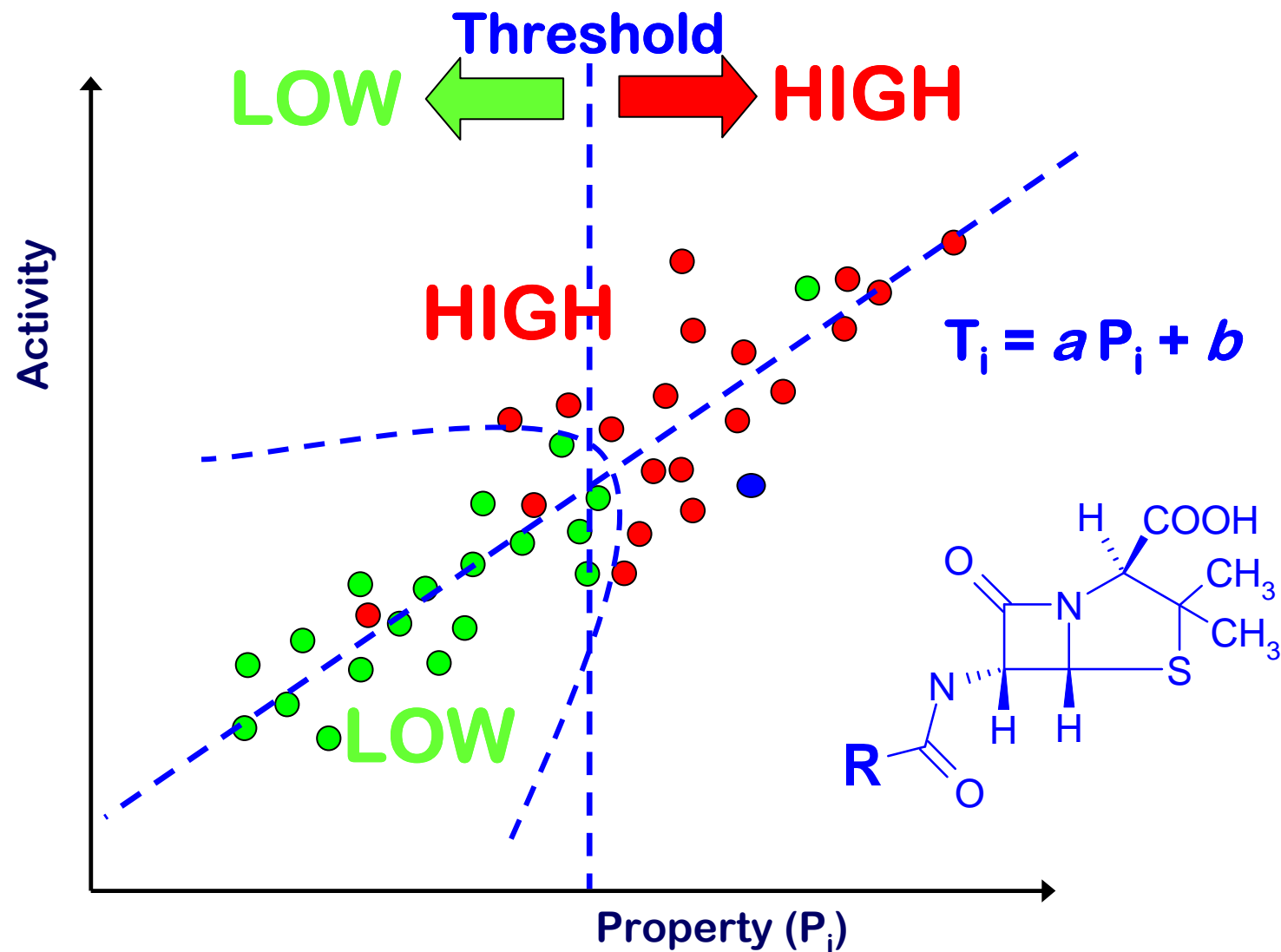
f(x)



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

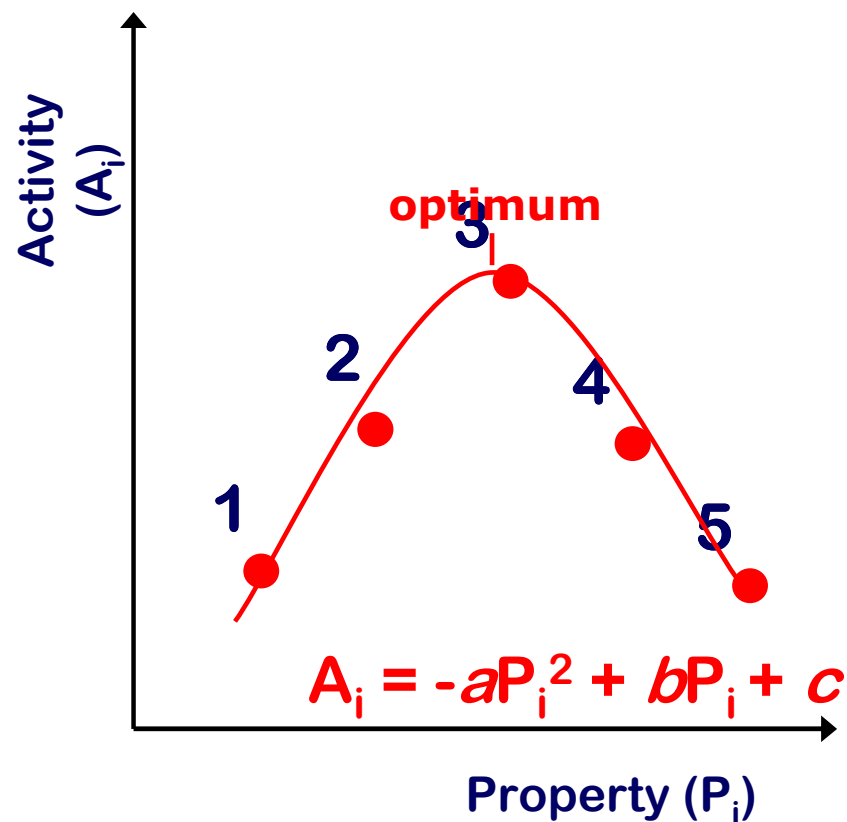
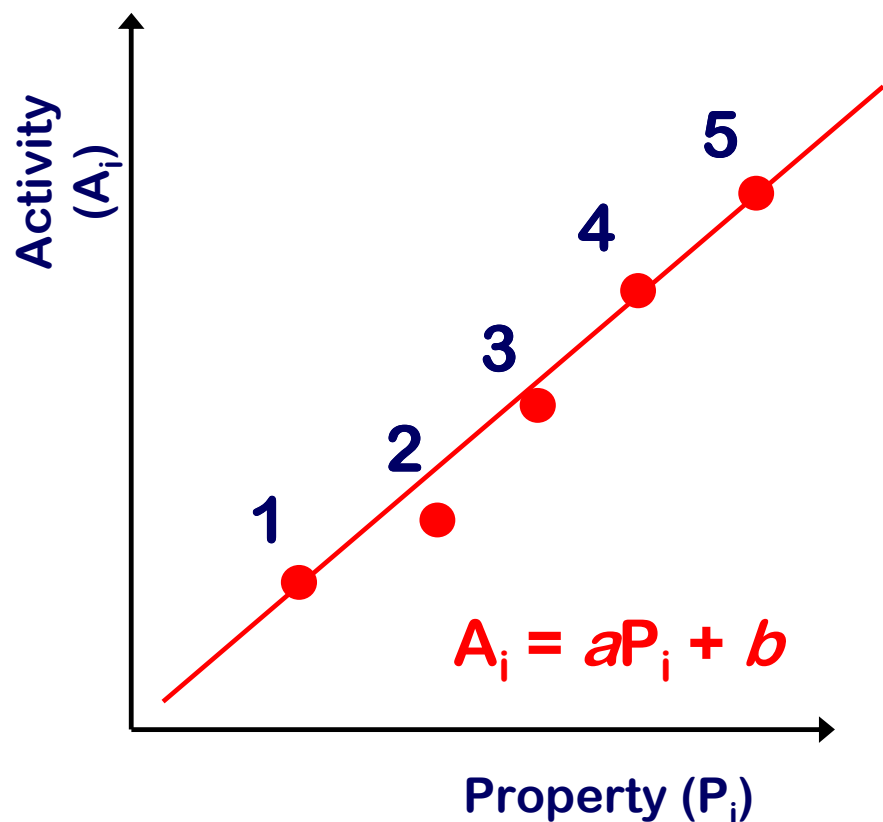


The scatter plot: the best place where explore (Q)SAR.



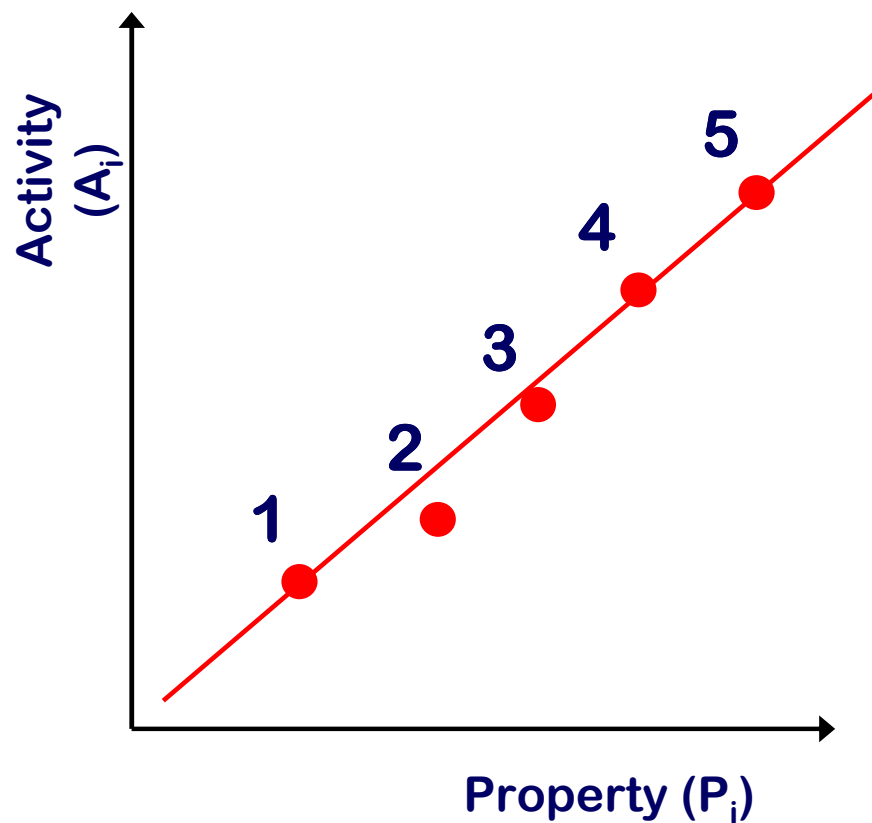


This is the base of a quantitative structure-activity relationship (QSAR): find *patterns*!





The beauty of mathematics:



$$A_i = a P_i + b$$

Discrete (few x - y correspondences)
Continuum (∞ x - y correspondences)

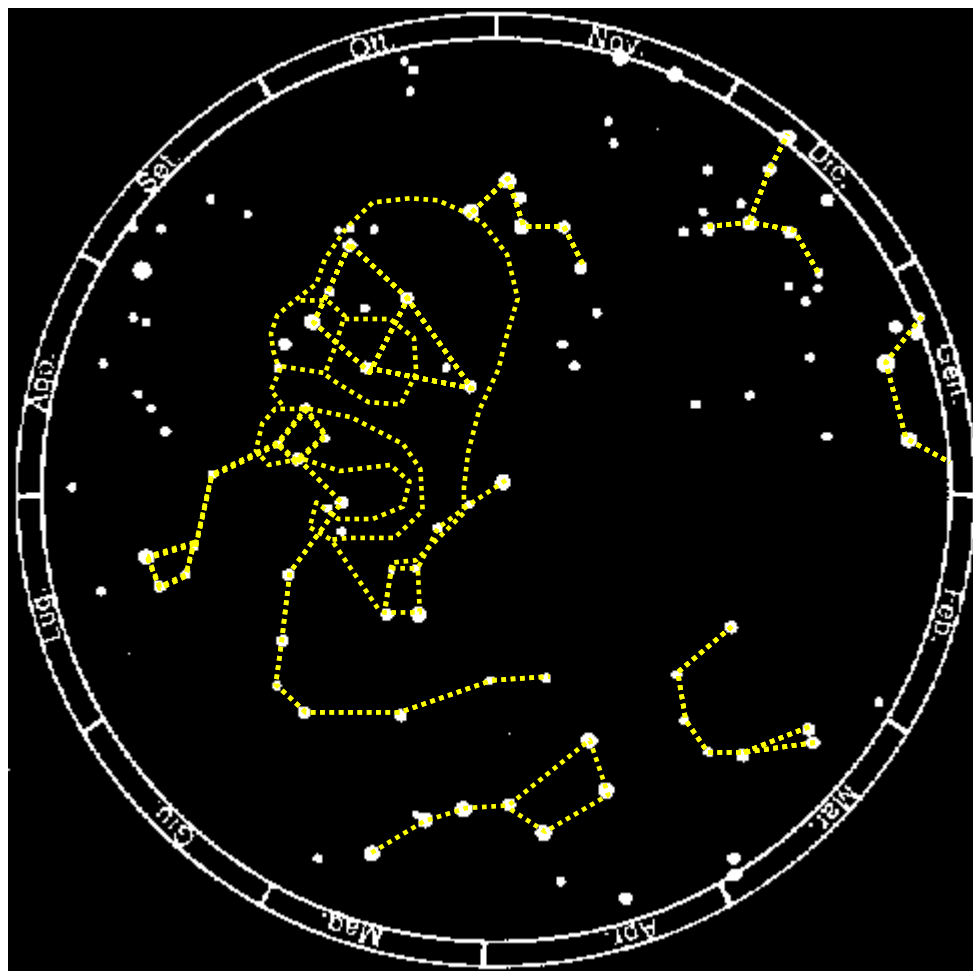


Patterns are beautiful:

- 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;



sometimes too beautiful...





The two statistical **gold** rules do build up linear models:

- For each independent variable (*molecular descriptor*) you need at least five (5) dependent variable values



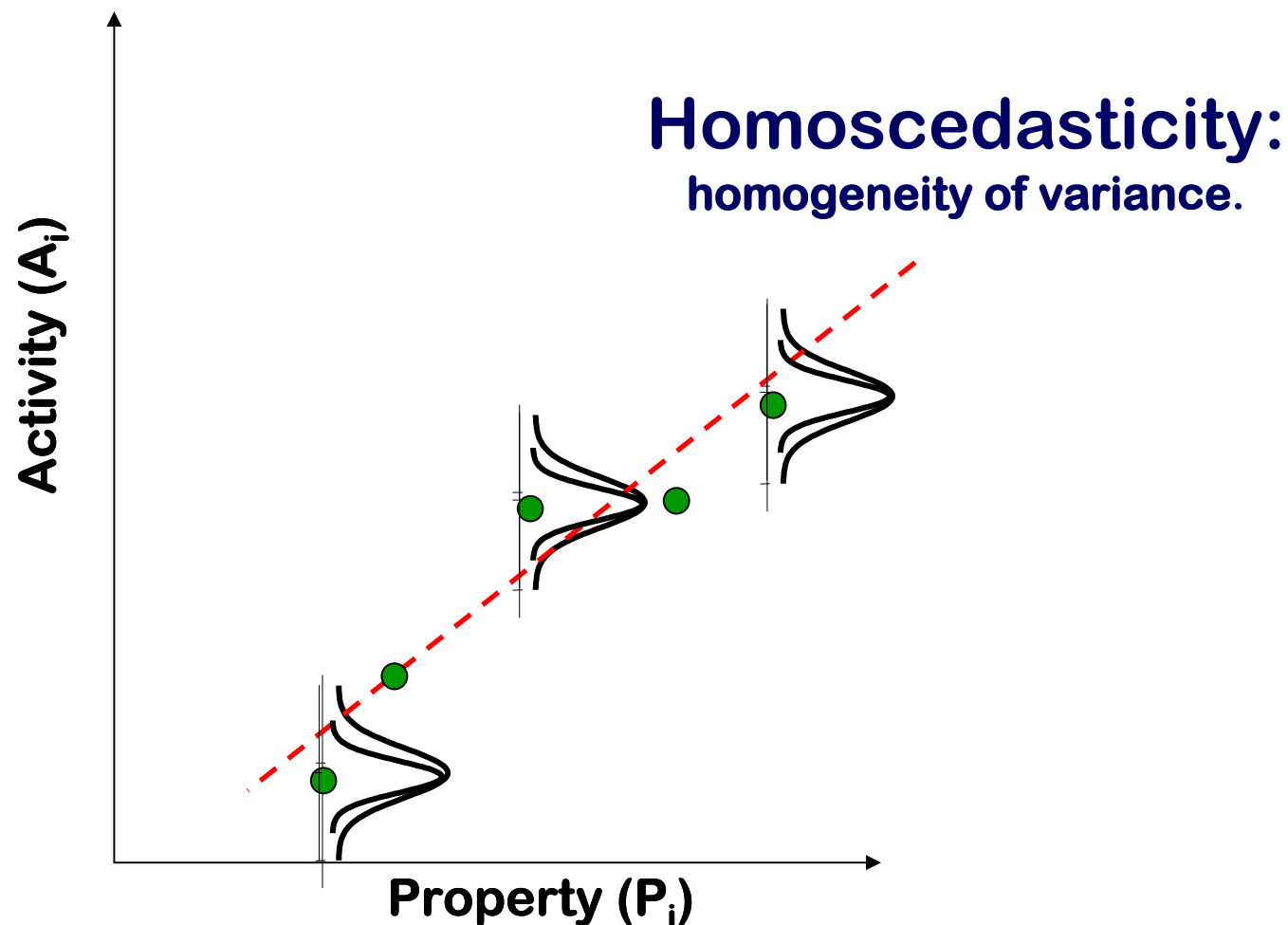
High accuracy, but low precision



High precision, but low accuracy



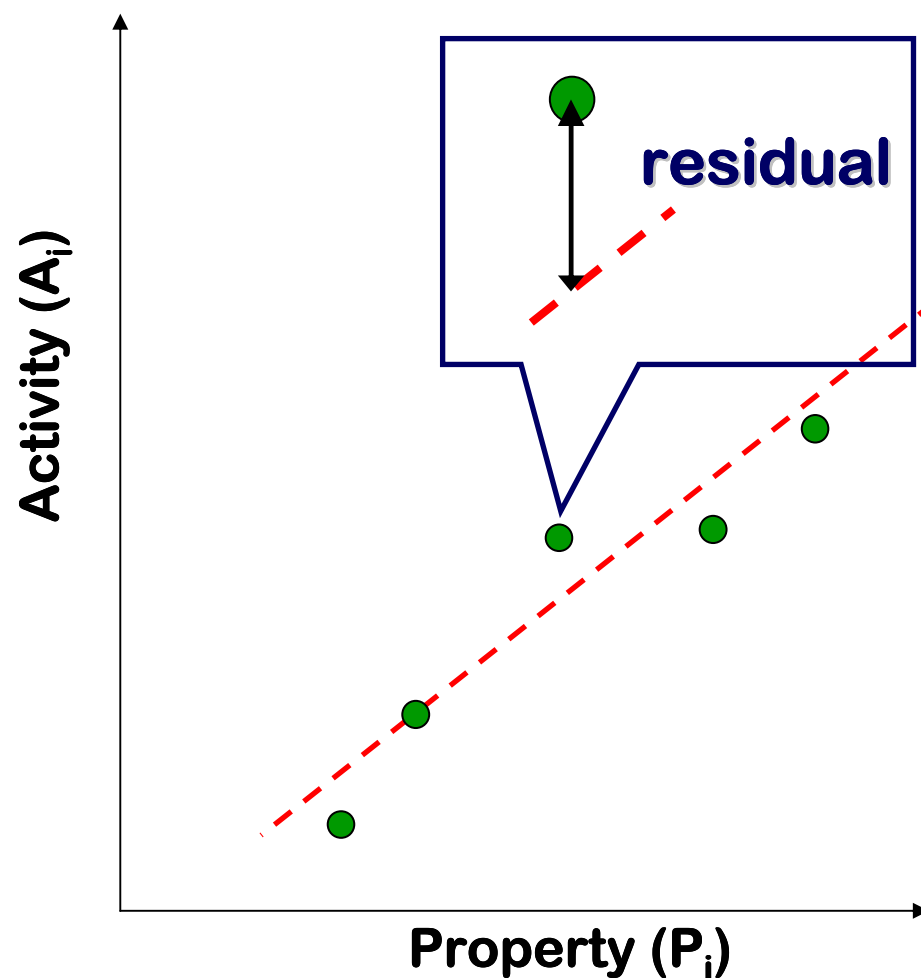
And how can we select the “*good*” linear model:



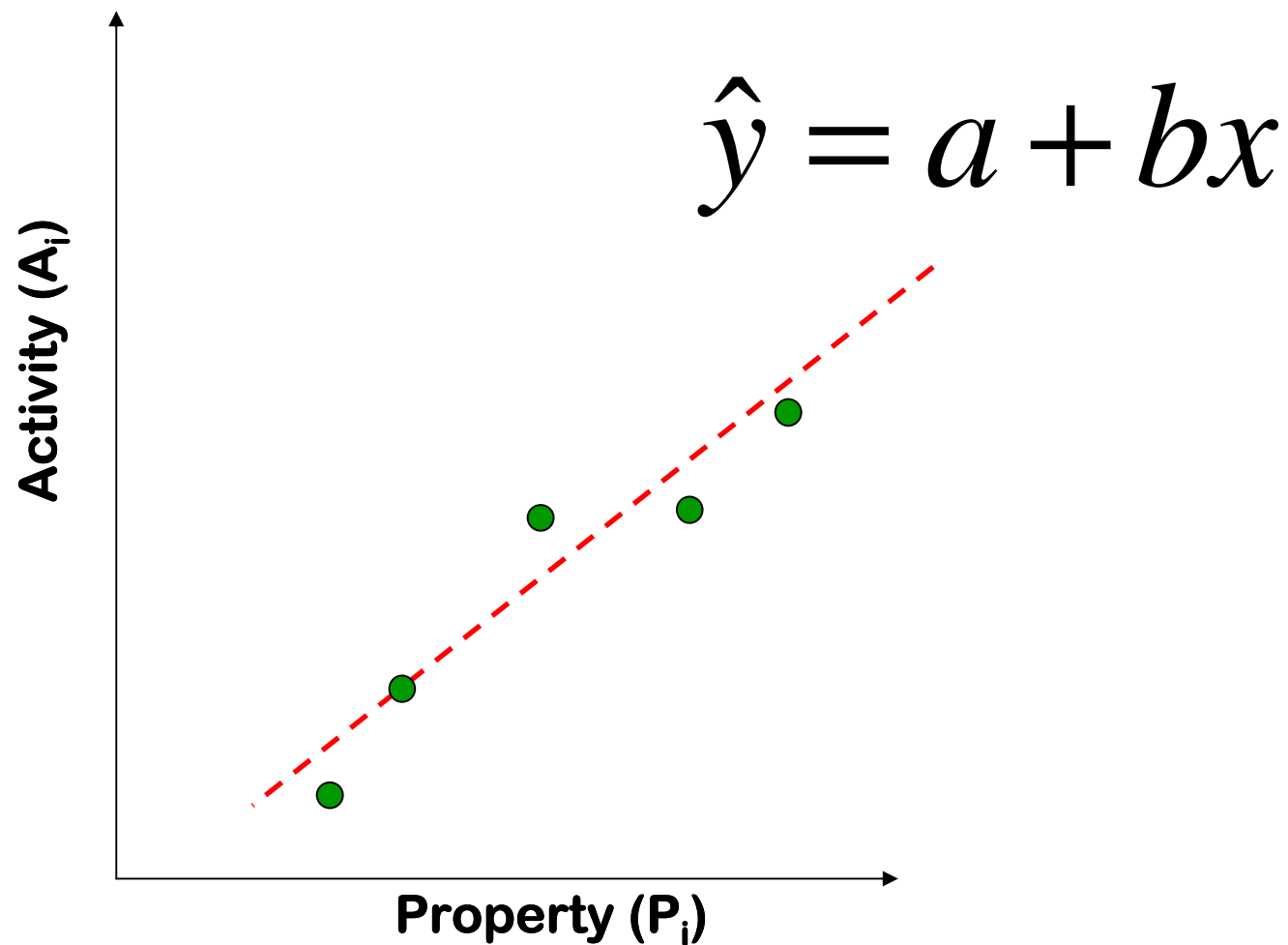


Do you remember... Least Squares Analysis?

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.

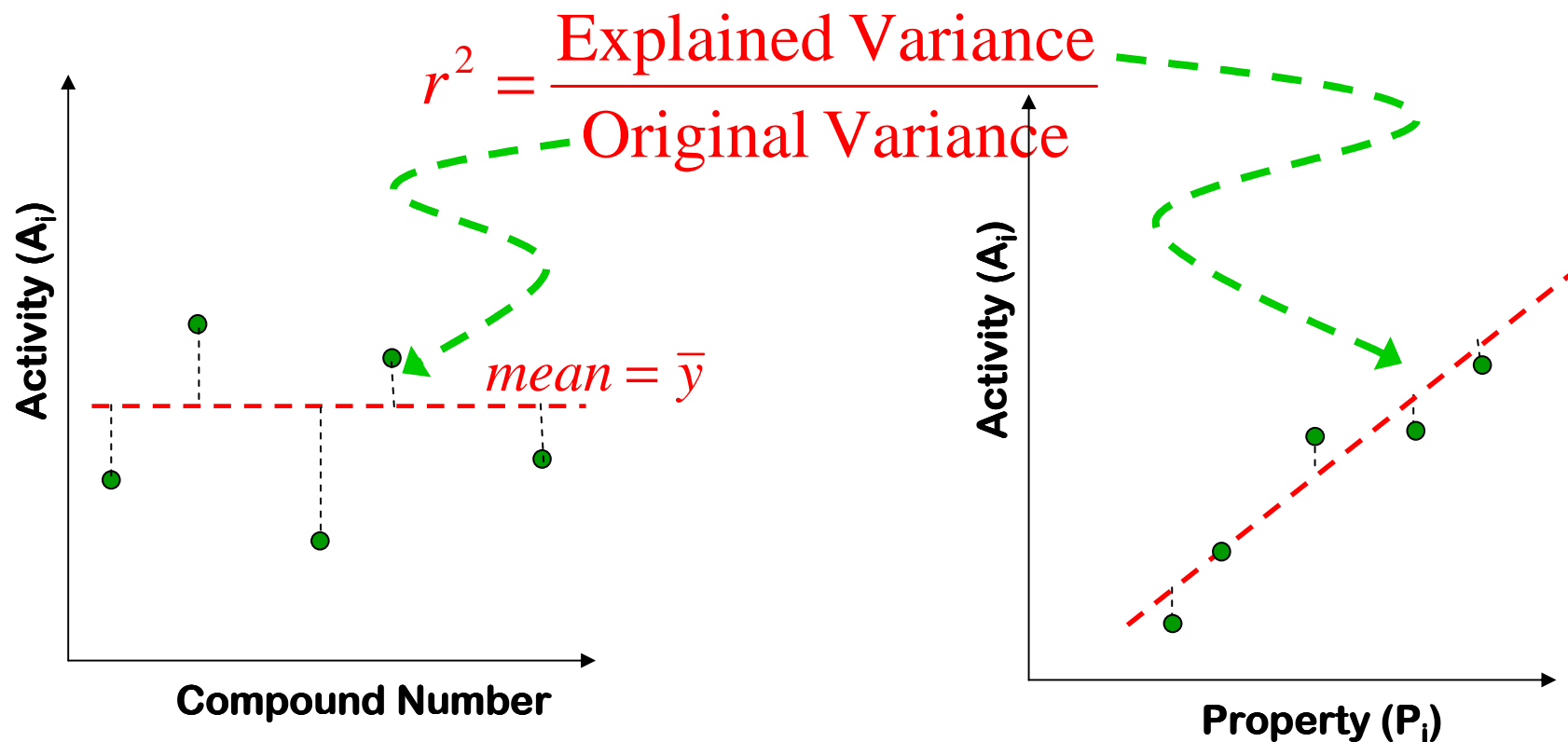


$$e = \hat{y} - y$$



$$b = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} \quad a = \frac{\sum y}{n} - b \frac{\sum x}{n}$$

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:



Original variance = Explained variance (*i.e.*, variance explained by the equation) + Unexplained variance (*i.e.*, residual variance around regression line)

Calculating r^2

- **Original variance:**

$$TSS = \sum_{i=1}^N (y_i - \bar{y})^2$$

- **Explained variance:**

$$ESS = \sum_{i=1}^N (y_{i,calc} - \bar{y})^2$$

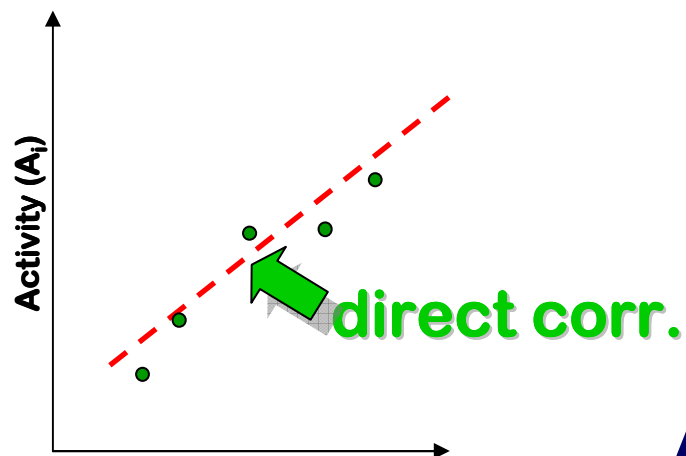
- **Variance around regression line:**

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

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

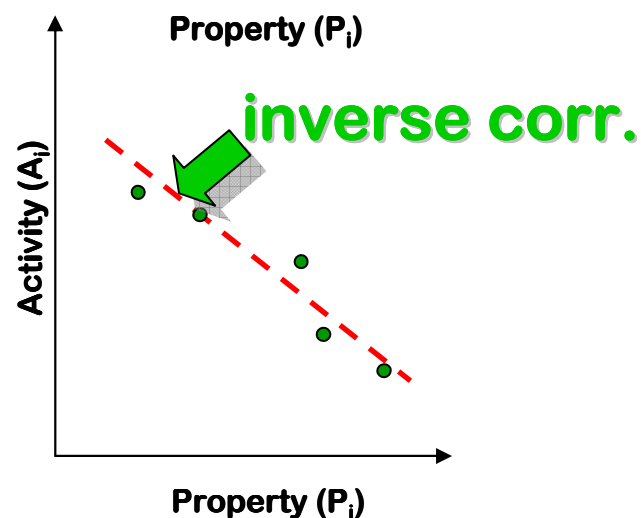
Possible values reported for r^2 fall between 0 and 1. For example: with r^2 of 0.83, you can say that 83% of the variability in activity can be explained by the different value of the selected molecular property. The remaining 17% of variability is due to other unexplained factors.

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:



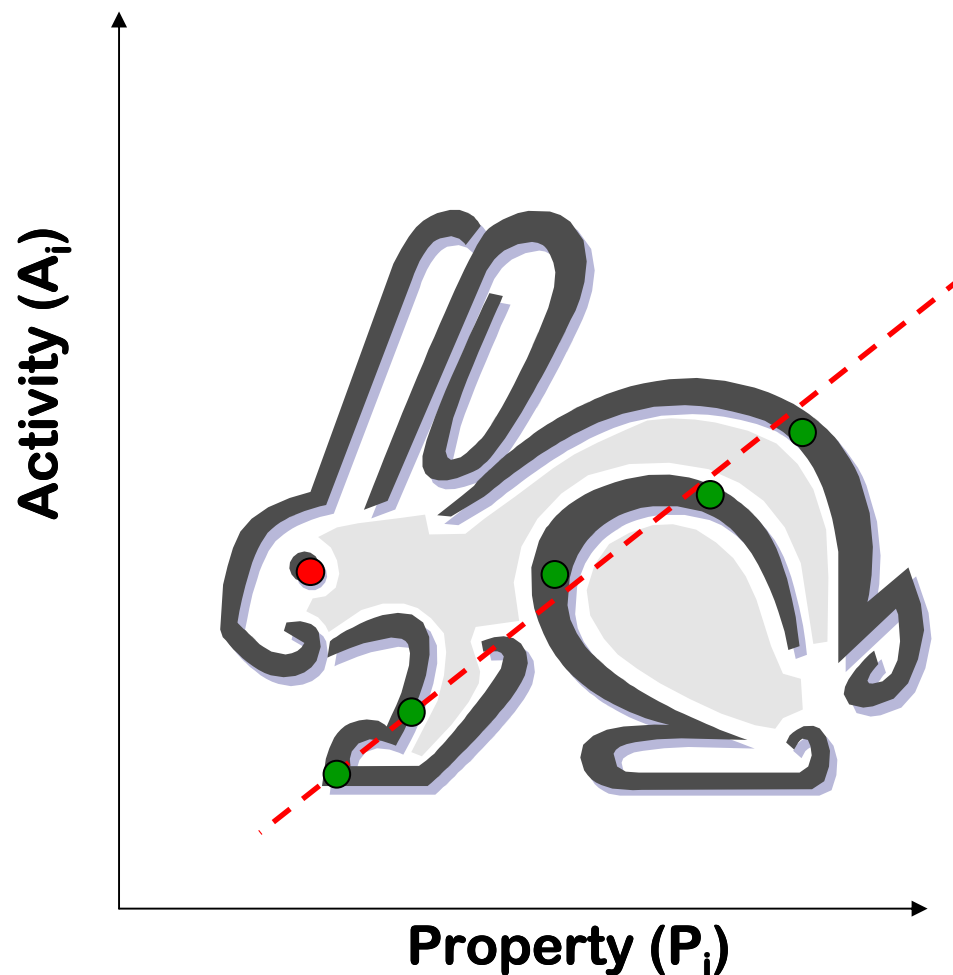
$$r^2 = \frac{ESS}{TSS} \quad r = \sqrt{\frac{ESS}{TSS}}$$

$$0 < r < 1$$



A perfect correlation of ± 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.

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

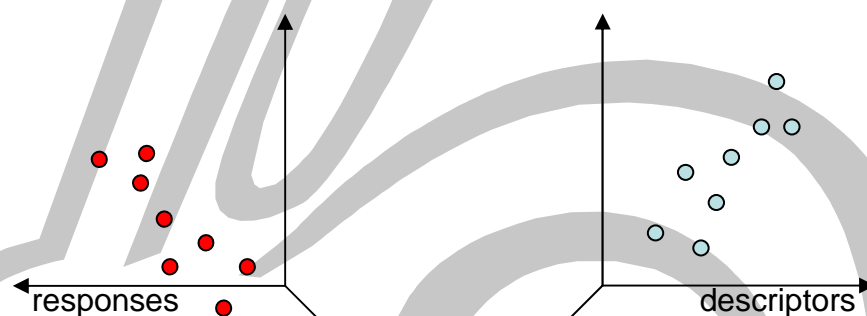


Be carefull... the rabbit is out there!!!



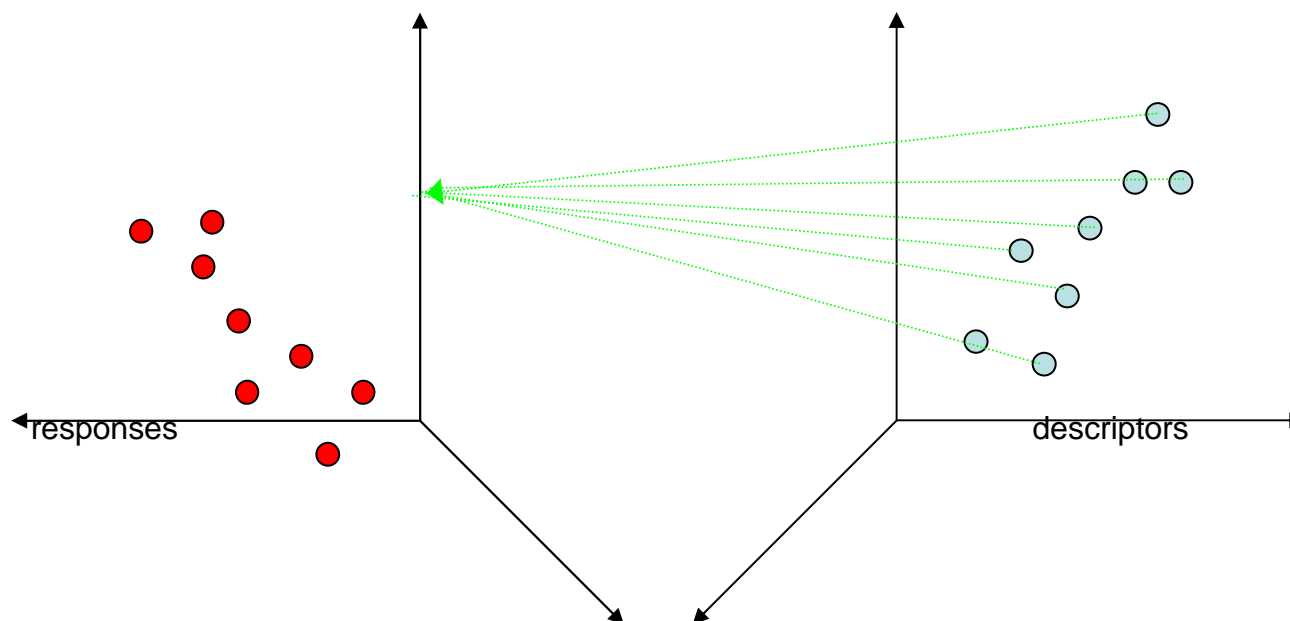
But usually in medchem, we have much more molecular descriptors than activity data. How can we statistically manage this situation?

Multiple Regression Analysis



$$\begin{aligned}y_1 &= a_1x_1 + b_1x_2 + \dots + m_1x_m + z_1 \\&\dots \\y_n &= a_nx_1 + b_nx_2 + \dots + m_nx_m + z_n\end{aligned}$$

Multiple Regression Analysis (MRA)



$$y_1 = a_1x_1 + b_1x_2 + \dots + m_1x_m + z_1$$

...

$$y_n = a_nx_n + b_nx_n + \dots + m_nx_m + z_n$$

Multiple Regression Analysis (MRA)

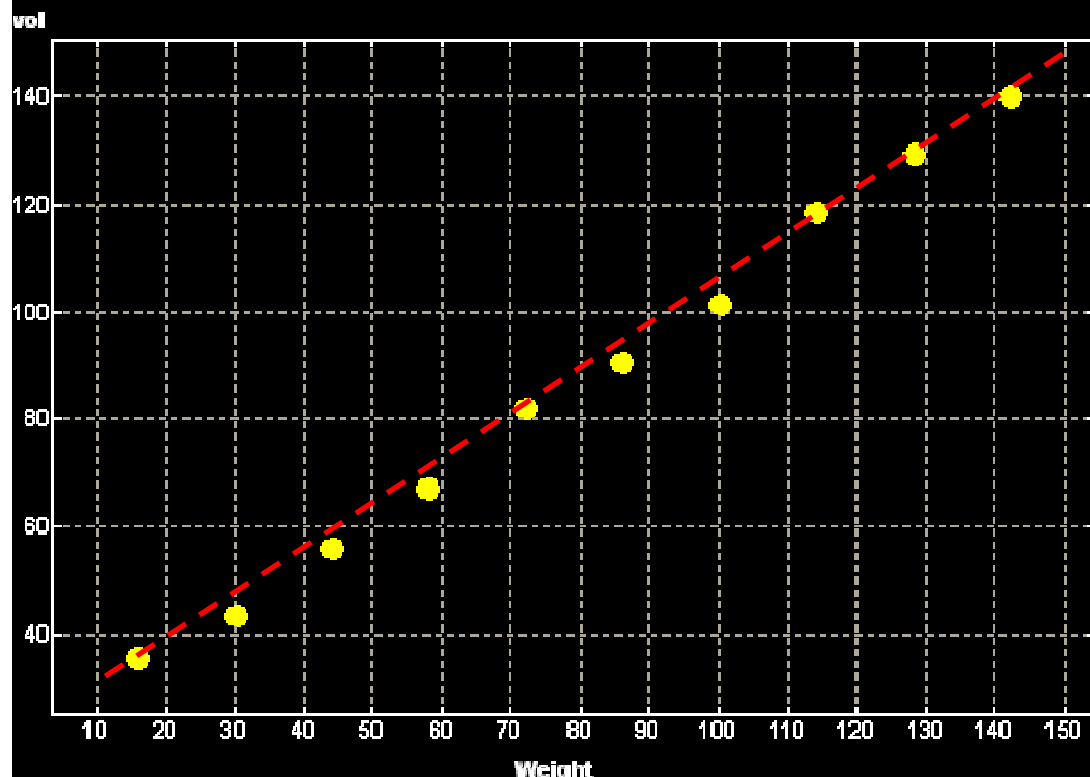
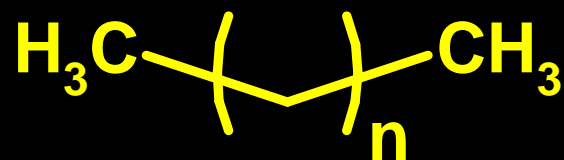
Requirements:

- There should be at least 5 times more samples than descriptors.
- Total number of descriptors should not exceed ~10 (looks the number of compounds you need!!!)
- Descriptors should be *uncorrelated*.



The second statistical **gold** rules do build up linear models:

- Having more the one molecular descriptors, the internal correlation (*cross-correlation*) between them has to be lower than 0.5

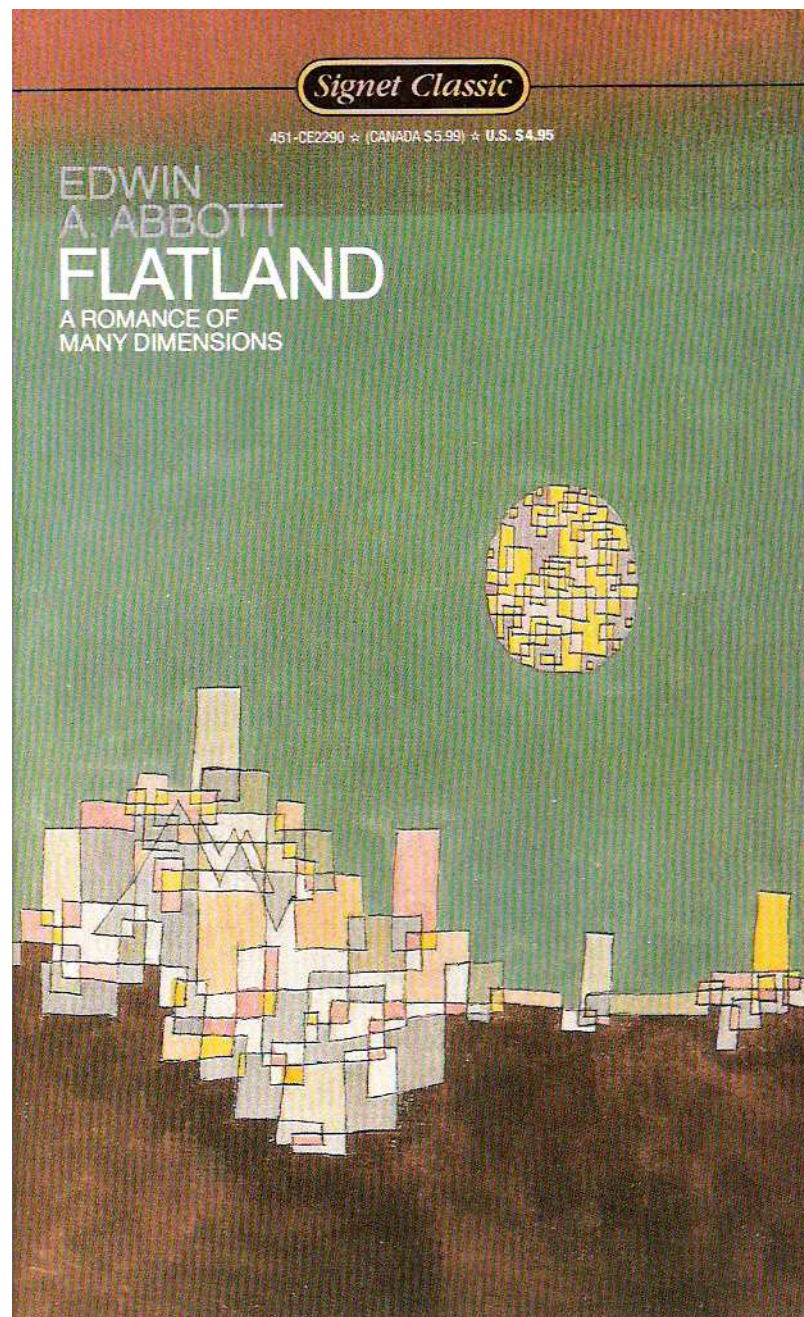


Descriptors:

- Molecular Volume
- Molecular Weight

$$r^2 = 0.9973$$

Considering this specific combination of dataset (aliphatic hydrocarbons and molecular descriptors) molecular volume and molecular weight are strongly correlate thus redundant!



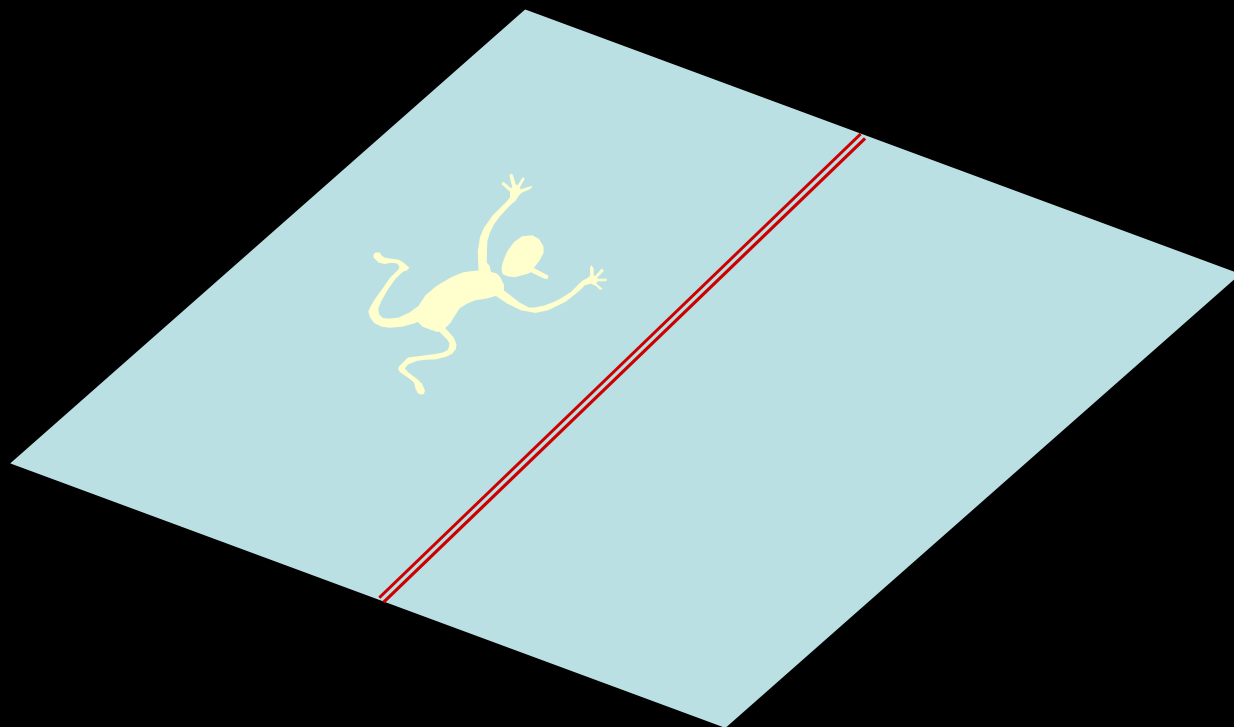
MS
M

Confidential and Property of ©2005 Molecular Modeling Section
Dept. Pharmaceutical and Pharmacological Sciences – University of Padova - Italy

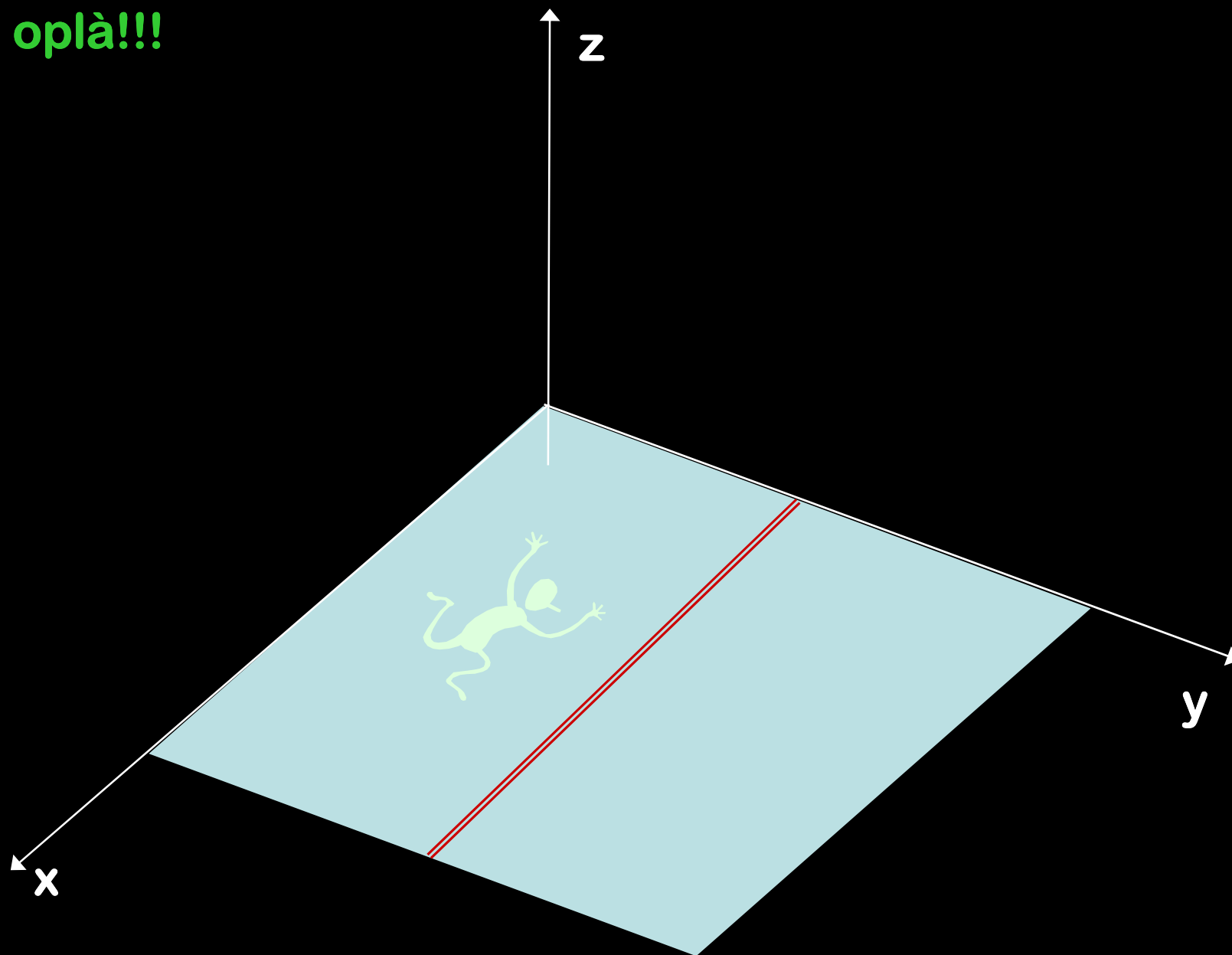
S. MORO – PSF – 2013/2014



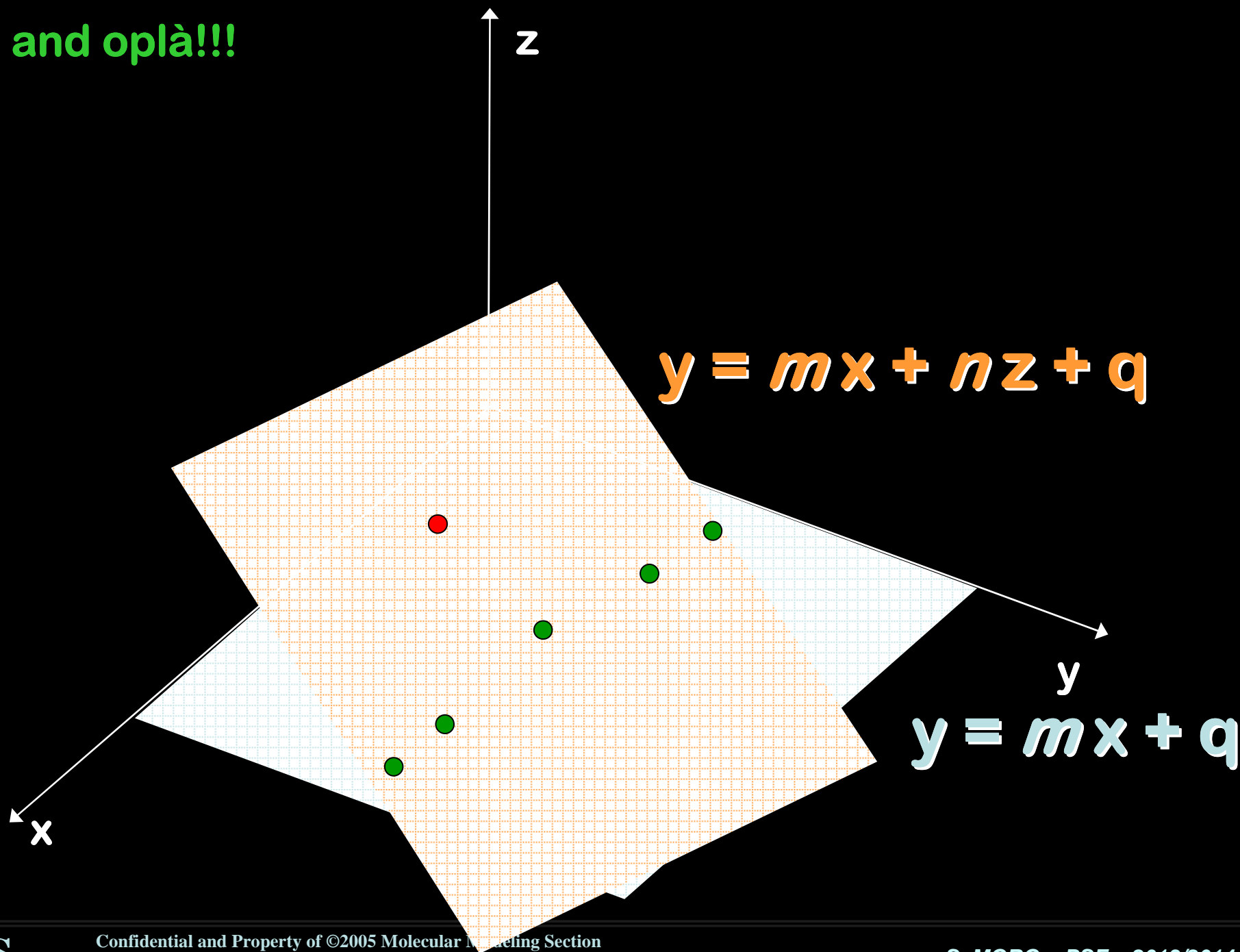
MRA approaches can transform the life in Flatland!



... opla!!!



... and oplà!!!



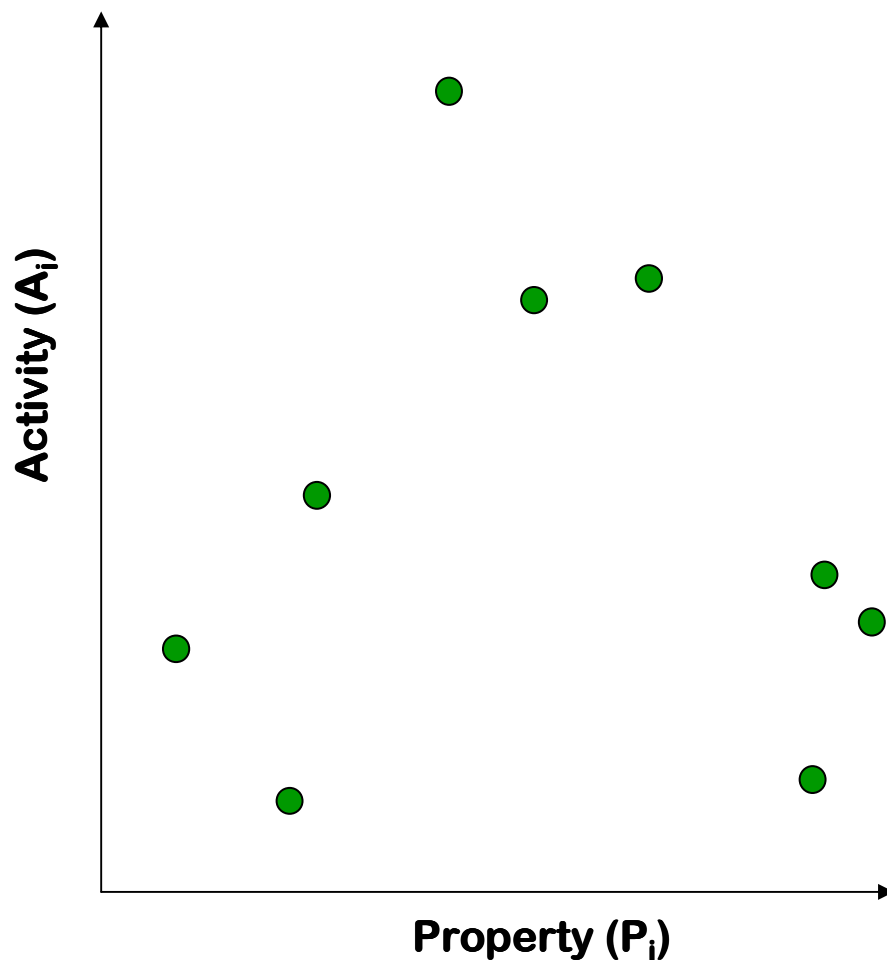


Here is the MRA nightmare:

r^2 will always increase as new descriptors are added.



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

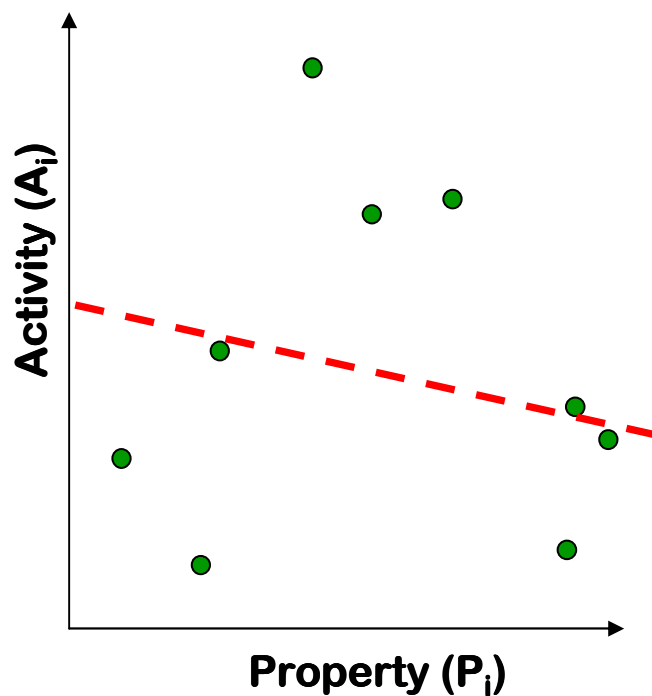


$$y = f(x) + \text{noise}$$

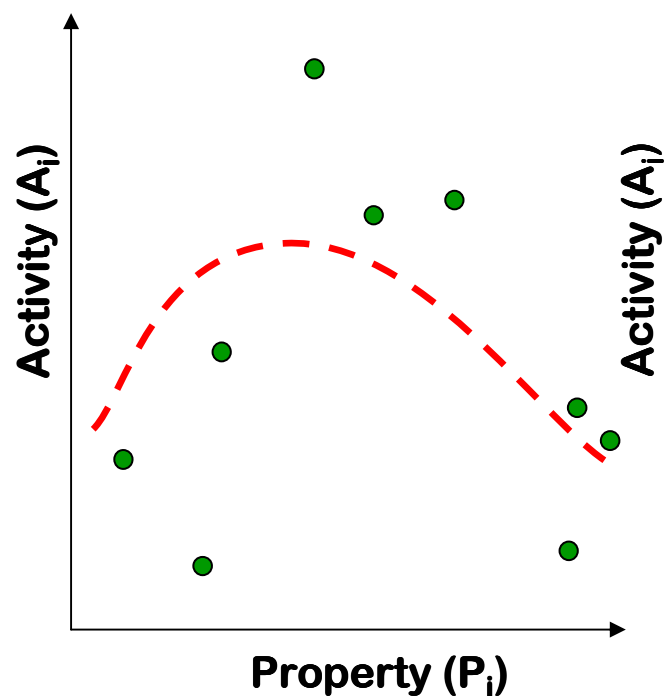
How we can deal with these data?

Let's consider three different methods...

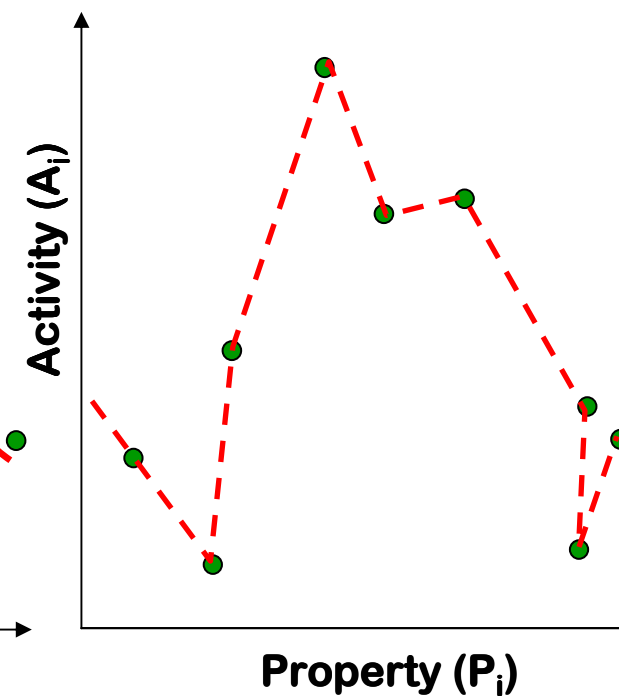
Linear



Quadratic



Joint-the-dots

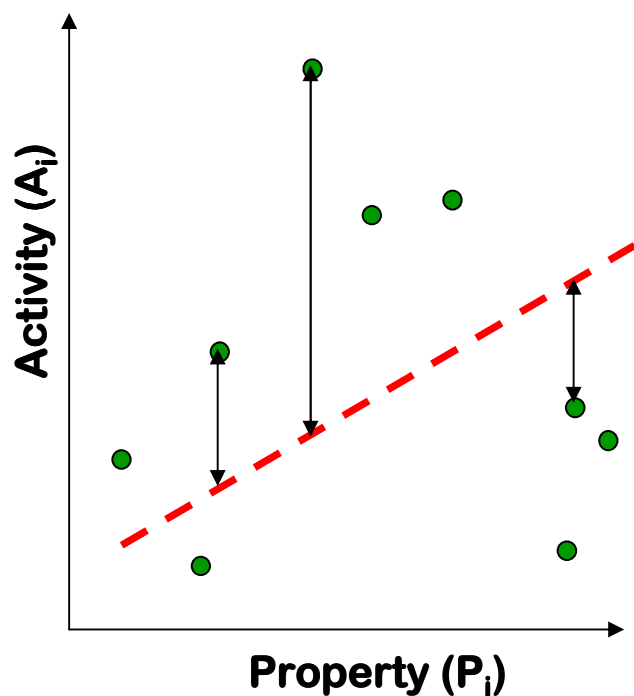


“How well are you going to predict future data drawn from the same distribution?”

Also known as
*piecewise linear non
parametric regression...*
if that makes you feel
better!!!



The “*test set*” method...looks this:



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

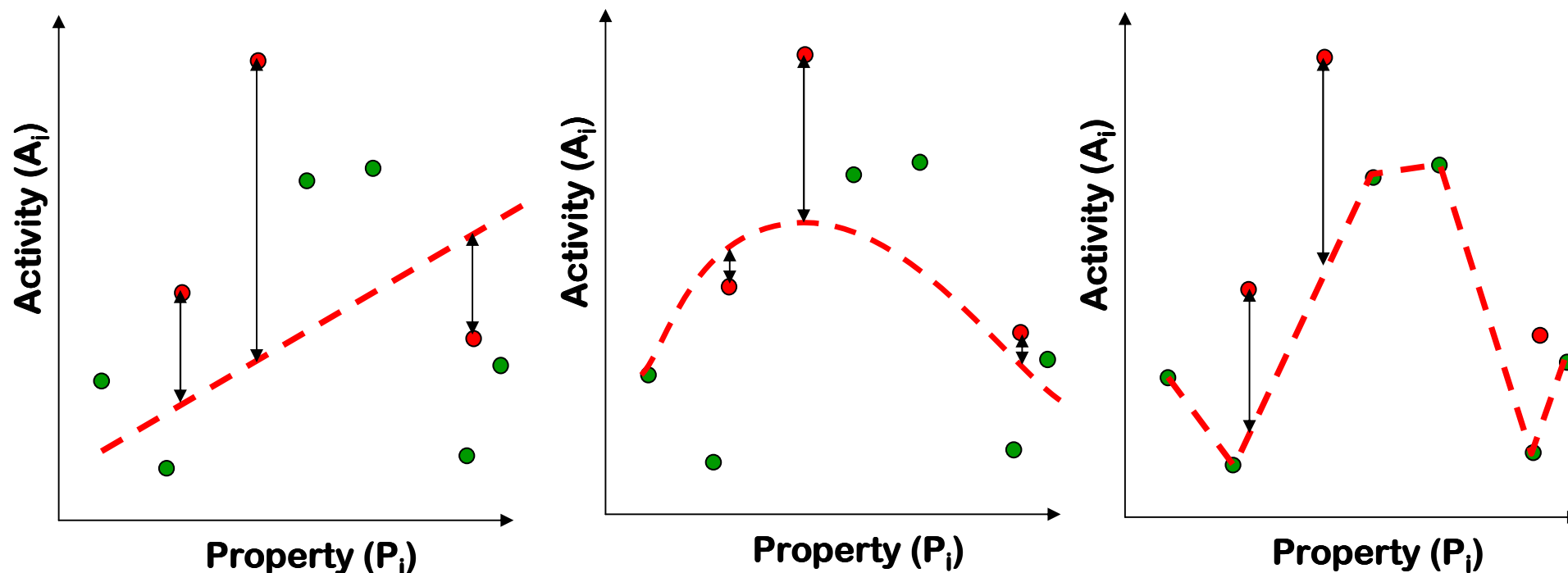
2. The remainder is a *training set*;

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

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

Mean Squared Error (MSE)

$$MSE = \frac{\sum (x_i - \bar{x})^2}{n}$$



MSE = 2.4

0.9

2.2

Good news:

- Very very simple;
- Can then simply choose the method with the best “*test set*” score.

Bad news:

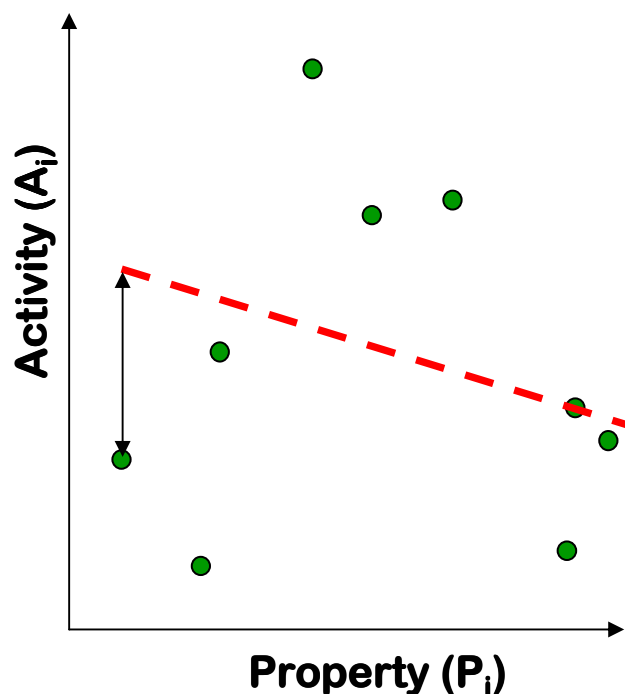
- Wastes data: we get an estimate of the best method to apply to 30% less data;
- If we don't have much data, our test-set might just be lucky or unlucky.



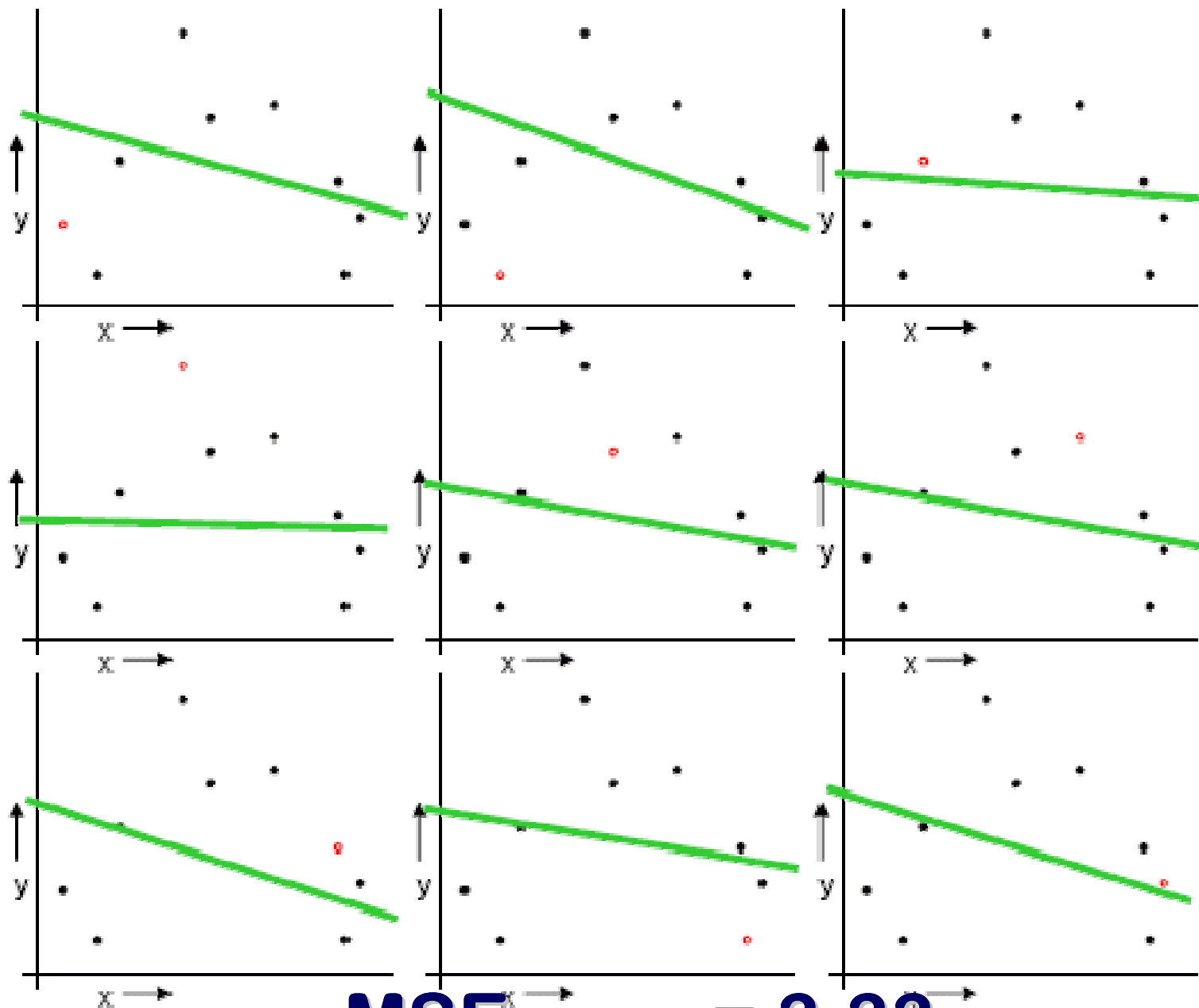
or “*LOOCV*” (Leave-One-Out Cross Validation) method... looks this:

For each data consider this loop:

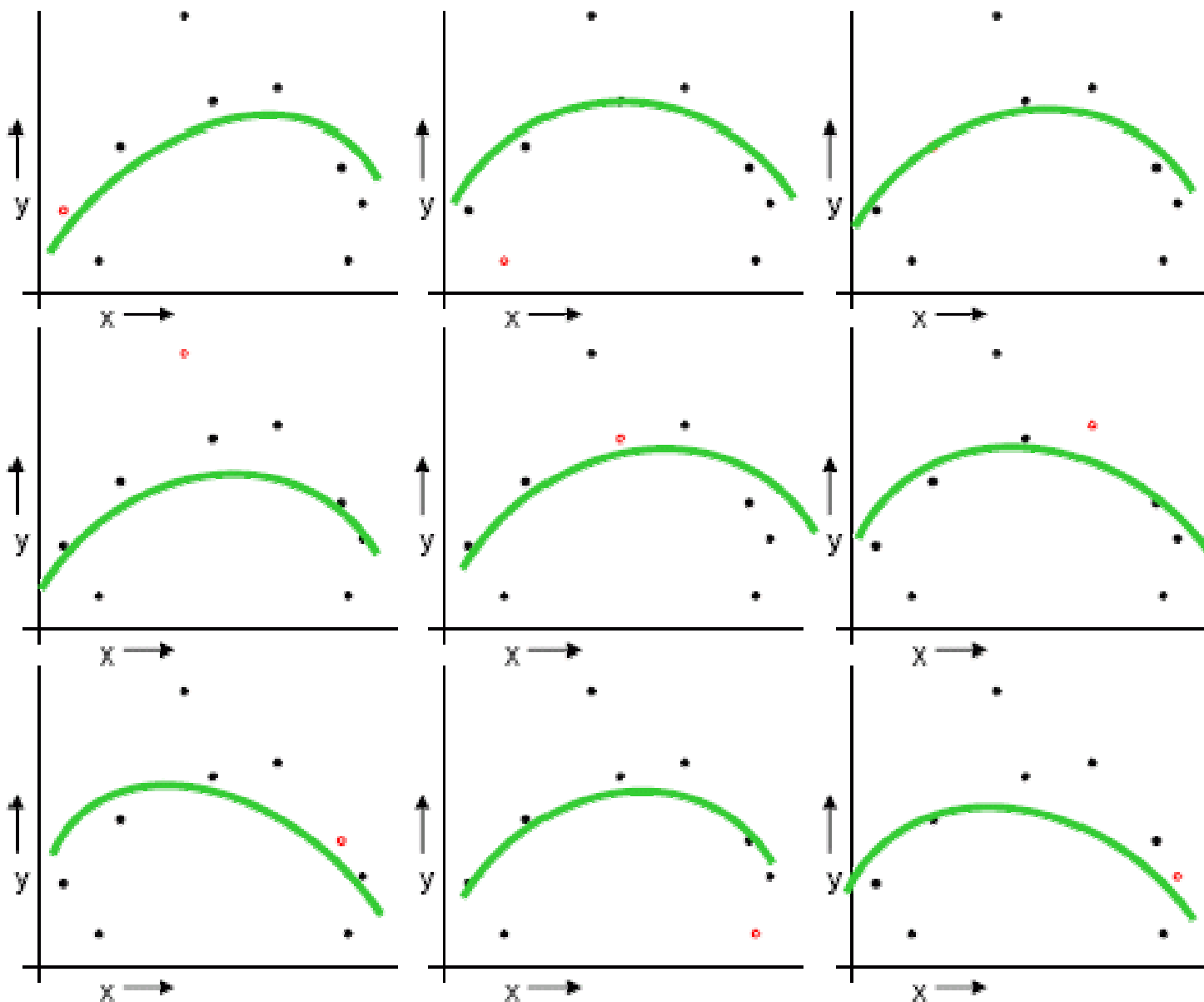
1. Select the **first** (x_i, y_i) data;
2. Temporary remove (x_i, y_i) from the data set;
3. Train on the remaining $n-1$ datapoints;
4. Note your error (x_i, y_i) ;



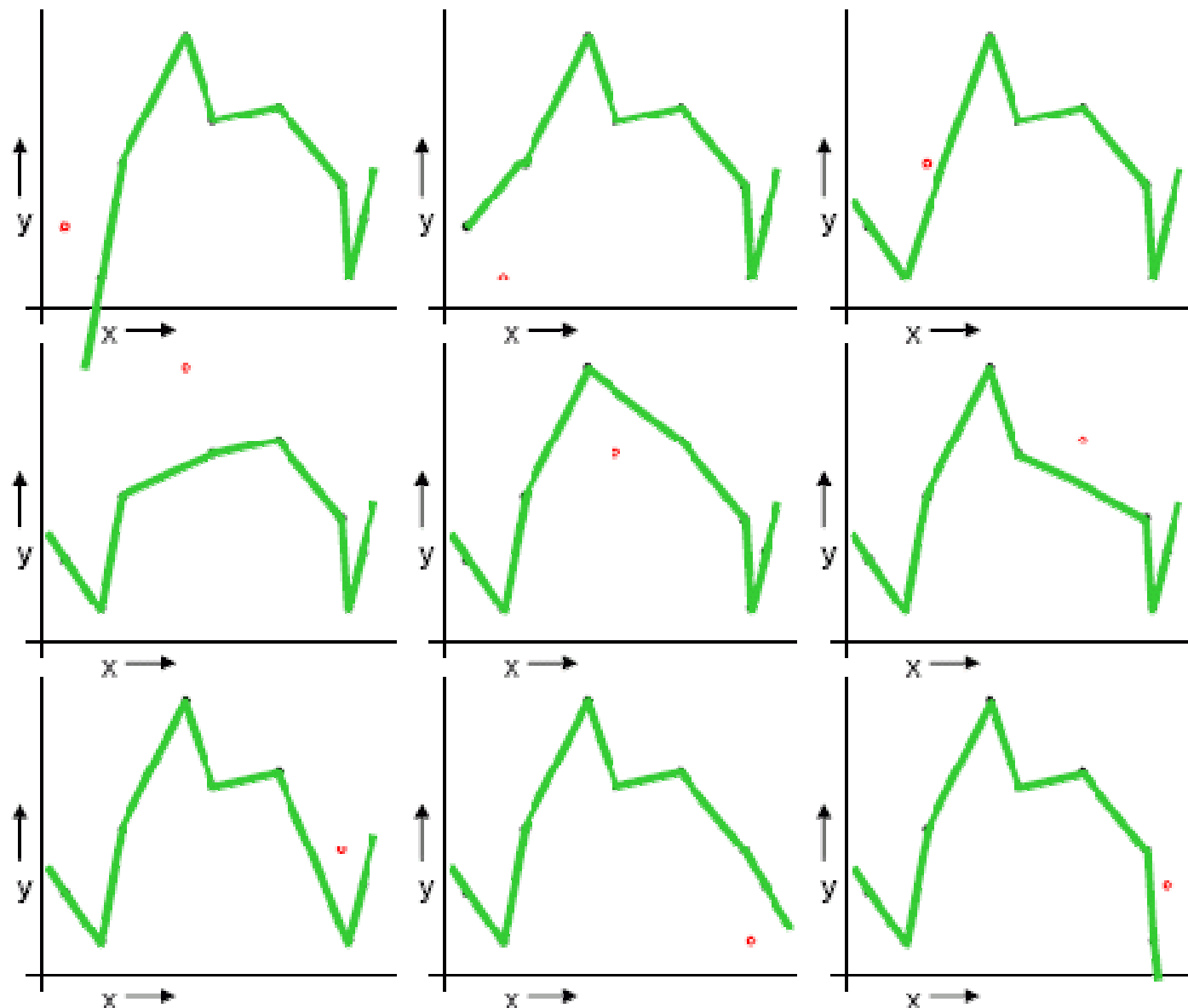
When you've done all points, report the mean squared errors (MSE).



$$\text{MSE}_{\text{Loocv}} = 3.33$$



$$\text{MSE}_{\text{Loocv}} = 0.96$$



$$\text{MSE}_{\text{Loocv}} = 2.12$$

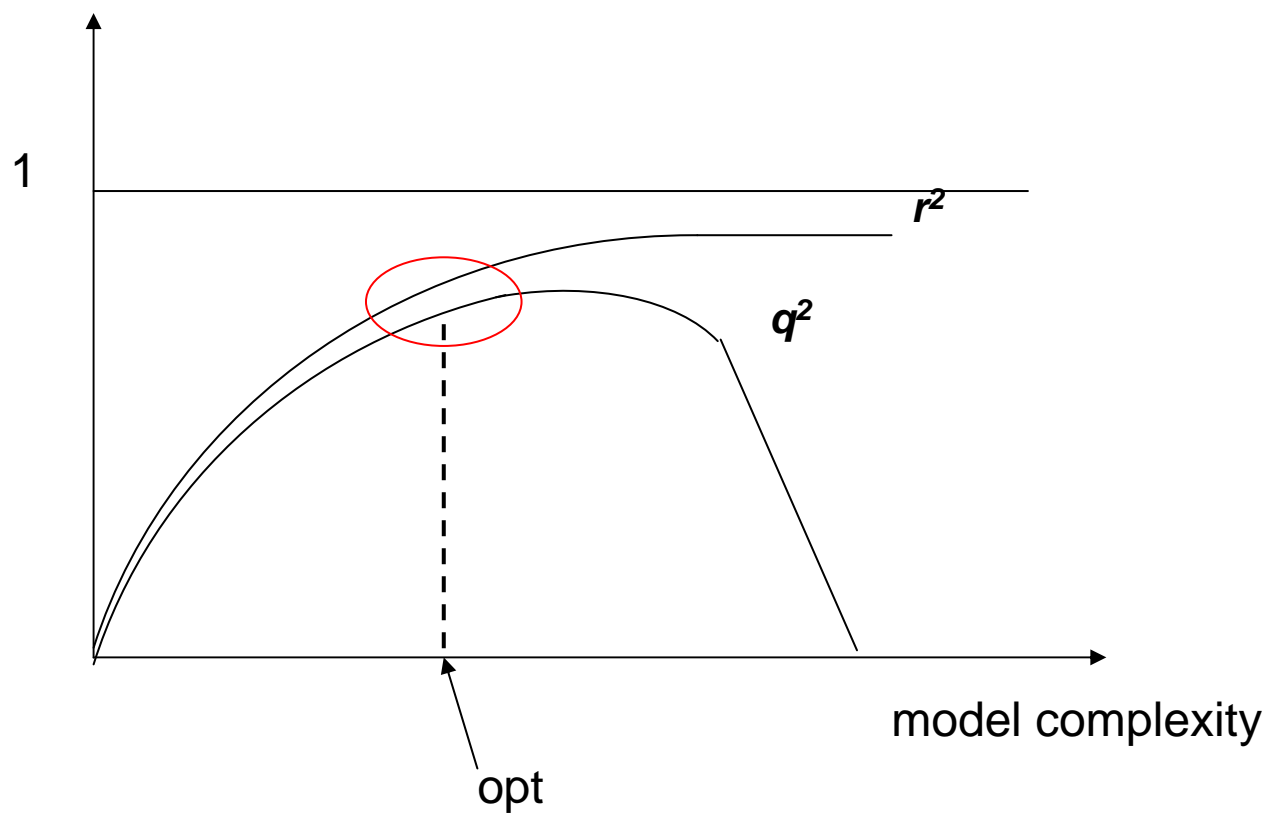


Cross-validation coefficient (Q^2)

$$Q^2 = 1 - \frac{PRESS}{\sum_{i=1}^N (y_i - \bar{y})^2}; \quad PRESS = \sum_{i=1}^N (y_{pred,i} - y_i)^2$$

$$r^2 = 1 - \frac{RSS}{\sum_{i=1}^N (y_i - \bar{y})^2}; \quad RSS = \sum_{i=1}^N (y_{calc,i} - y_i)^2$$

Q^2 initially increases as more parameters are added but then starts to decrease indicating data over fitting. Thus Q^2 is a better indicator of the model quality.





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



Another important consideration using MRA technique:

#	MR	logP	Volume	PM	Surface	density	n. X atoms
CH ₂ Cl ₂	1,62959	1,30436	67,1359	84,933	176,1169	1,63677	3
CHCl ₃	2,01731	1,73808	75,7514	119,378	186,8237	2,03576	4
CCl ₄	2,35508	2,42116	83,2702	153,823	194,0565	2,3224	5
CF ₃ CHBrCl	2,35642	2,36112	88,098	197,381	206,6438	2,85284	7
CHCl ₂ CHCl	2,92829	2,49472	94,2376	167,85	215,3294	2,26061	6
Cl ₂ C=CHCl	2,46705	2,28836	115,831	131,389	241,6985	1,42863	5
CCl ₂ =CCl ₂	2,82835	3,37472	132,106	165,834	257,2367	1,46129	6

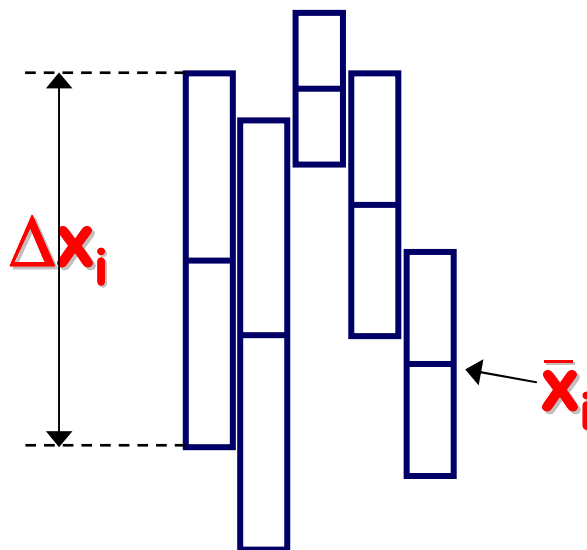
$$\bar{x}_i = \quad 2.37 \quad 2.28 \quad 93.77 \quad 145.80 \quad 211.13 \quad 1.99 \quad 5.10$$

$$\Delta x_i = \quad 1.30 \quad 2.07 \quad 44.00 \quad 112.45 \quad 81.12 \quad 1.42 \quad 4.00$$

data scaling and data centering

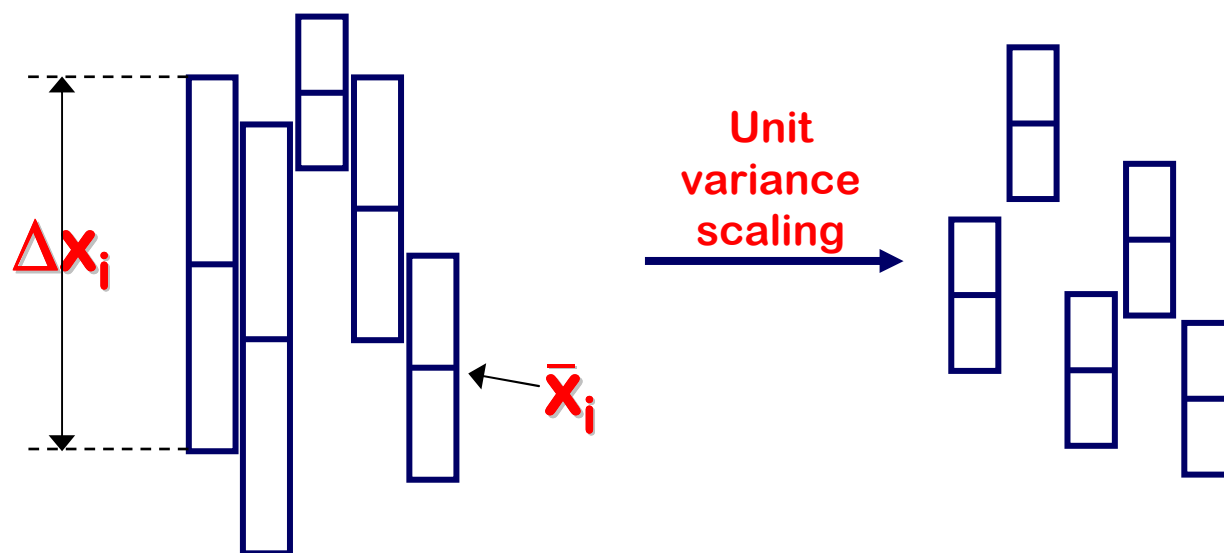
data scaling and data centering

- Each independent variable influences the model according to its variance.
- Thus scaling corresponds to the assumption that all variables are *a priori* equally important.



data scaling and data centering

- **Unit variance scaling:** multiply each column by $1/\sigma_i$, σ_i being the standard deviation.



$$\sigma_i = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} \quad \text{where } n \text{ is the number of data taken}$$



Back to the real case:

#	MR	logP	Volume	PM	Surface	density	n. X atoms
CH ₂ Cl ₂	1,62959	1,30436	67,1359	84,933	176,1169	1,63677	3
CHCl ₃	2,01731	1,73808	75,7514	119,378	186,8237	2,03576	4
CCl ₄	2,35508	2,42116	83,2702	153,823	194,0565	2,3224	5
CF ₃ CHBrCl	2,35642	2,36112	88,098	197,381	206,6438	2,85284	7
CHCl ₂ CHCl	2,92829	2,49472	94,2376	167,85	215,3294	2,26061	6
Cl ₂ C=CHCl	2,46705	2,28836	115,831	131,389	241,6985	1,42863	5
CCl ₂ =CCl ₂	2,82835	3,37472	132,106	165,834	257,2367	1,46129	6

$$\bar{x}_i = \quad 2.37 \quad 2.28 \quad 93.77 \quad 145.80 \quad 211.13 \quad 1.99 \quad 5.10$$

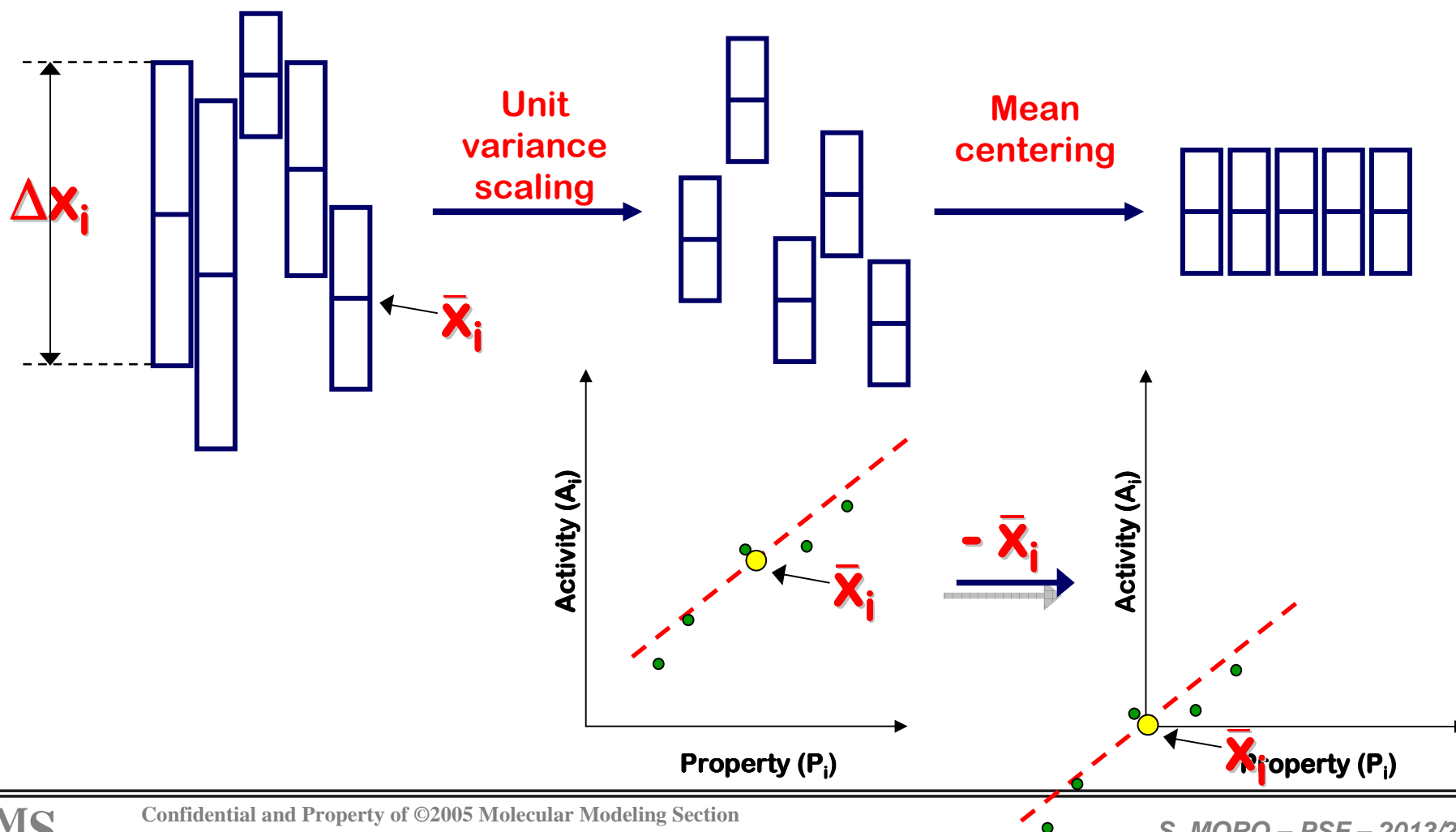
$$\Delta x_i = \quad 1.30 \quad 2.07 \quad 44.00 \quad 112.45 \quad 81.12 \quad 1.42 \quad 4.00$$

$$\sigma_i = \quad 0.45 \quad 0.65 \quad 22.85 \quad 37.02 \quad 29.46 \quad 0.52 \quad 1.34$$

$$\Delta x_i / \sigma_i = \quad 2.89 \quad 3.18 \quad 1.92 \quad 3.04 \quad 2.75 \quad 2.73 \quad 2.98$$

data scaling and data centering

- **Mean centering:** subtract from each column its average value.

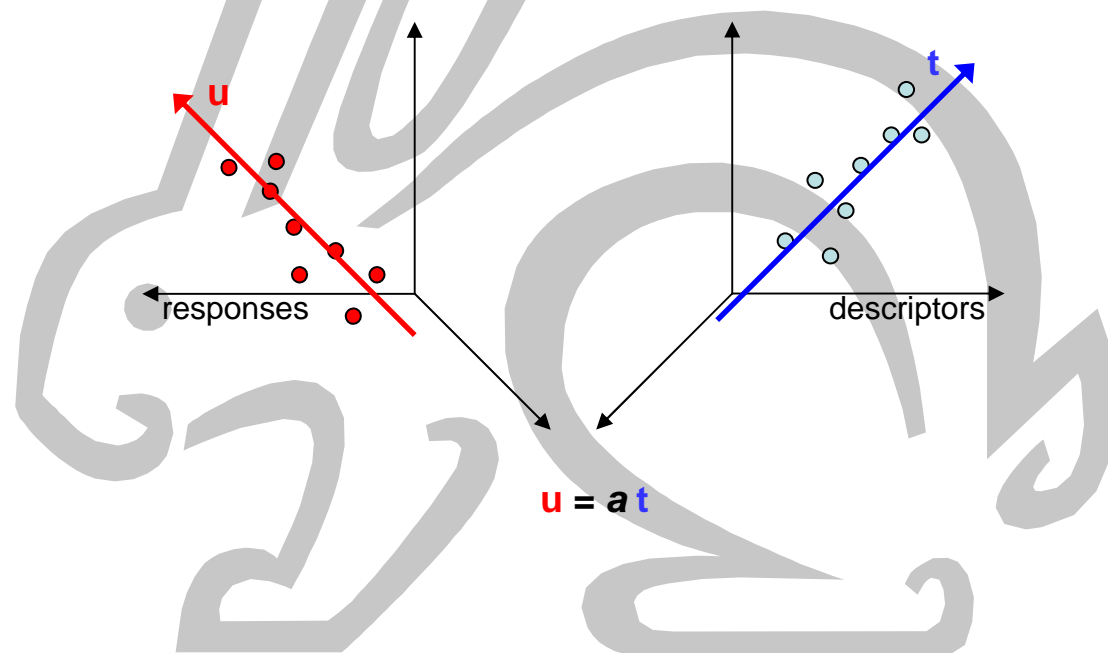




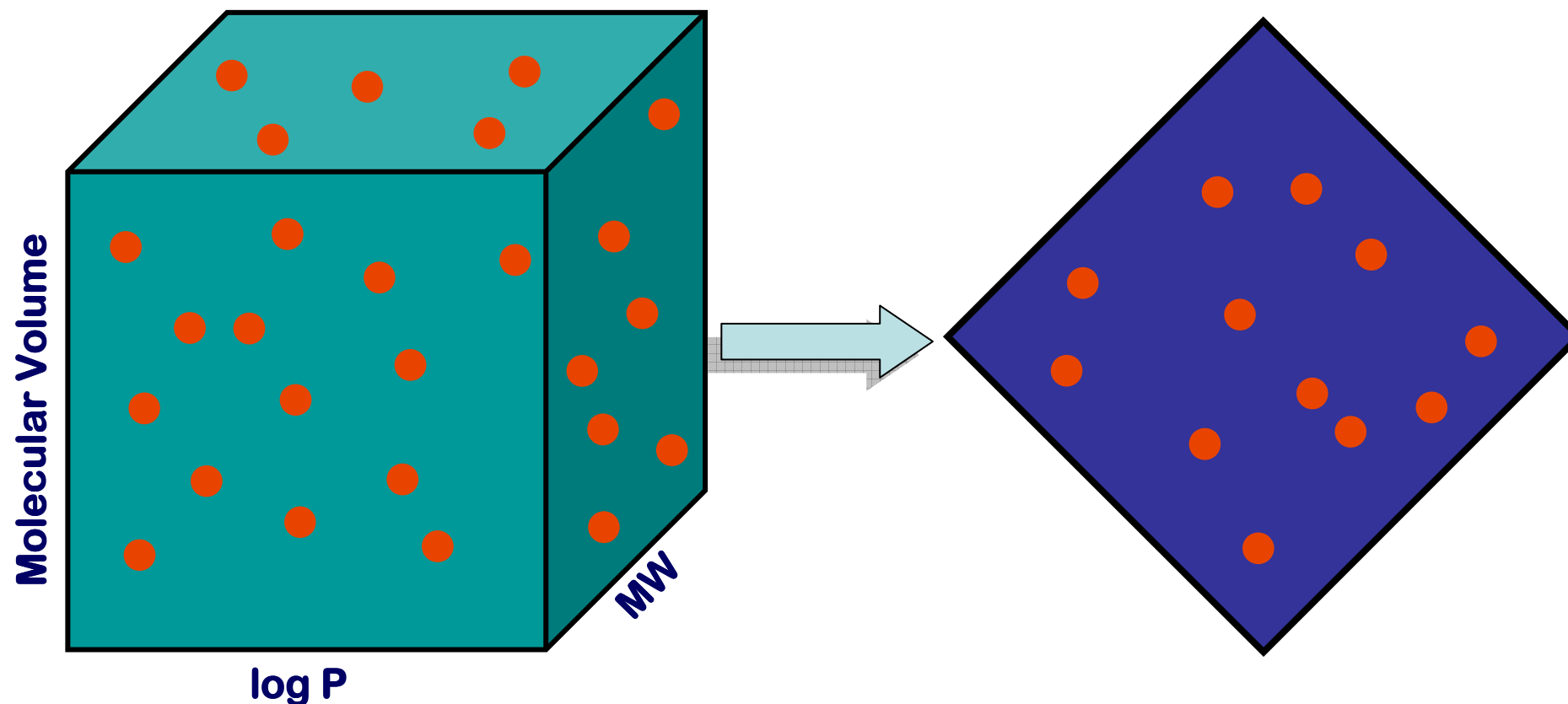
MRA should be a suitable tool only if these criteria are respected:

1. Good ratio between independent and dependent variables;
2. Statistical significance of the regression coefficient;
3. The magnitude of the typical effect " $b_i x_i$ ";
4. Any cross-correlation with other terms.

Principal Component Regression



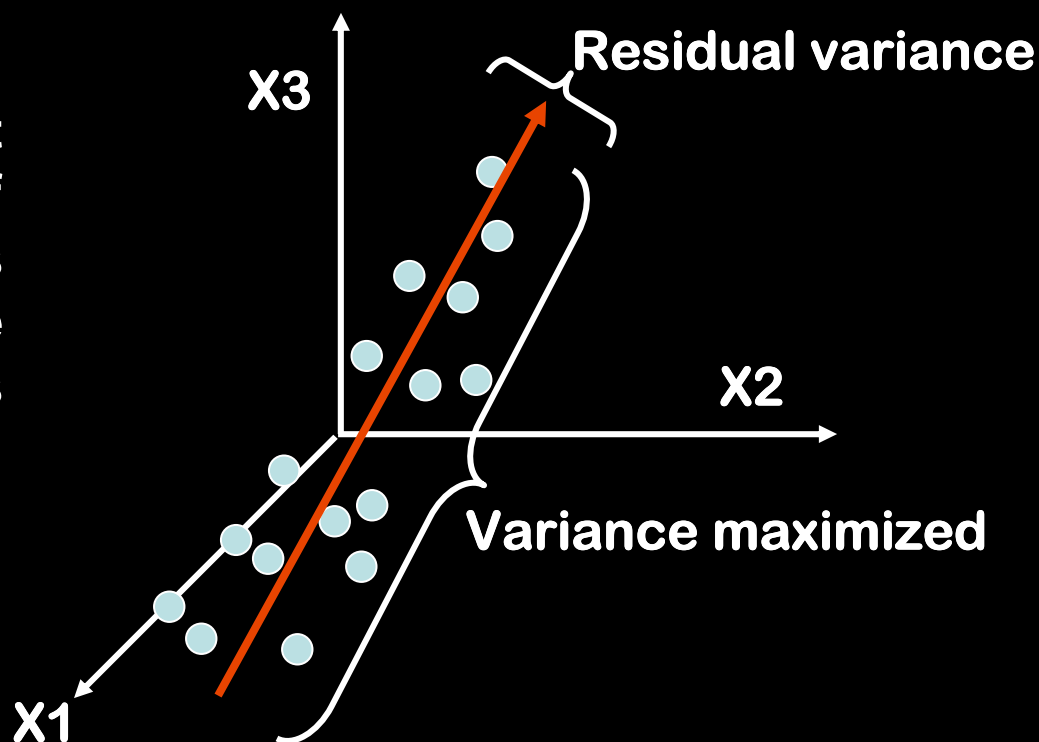
Data Presentation: Property Space



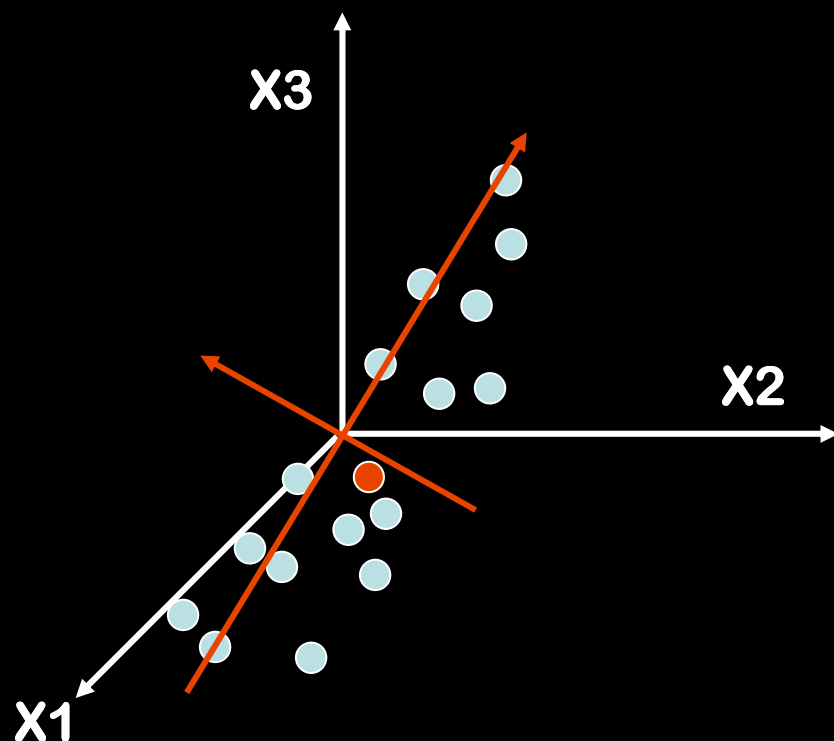
Principle Component Analysis (PCA)

PCA finds lines, planes and hyperplanes in the originally K-dimensions space that approximate the data as well as possible in the least square sense. In such a case, the variance in the original data is maximized.

A line that is the least squares approximation of a set of data points makes the variance of the coordinates on the line as large as possible.



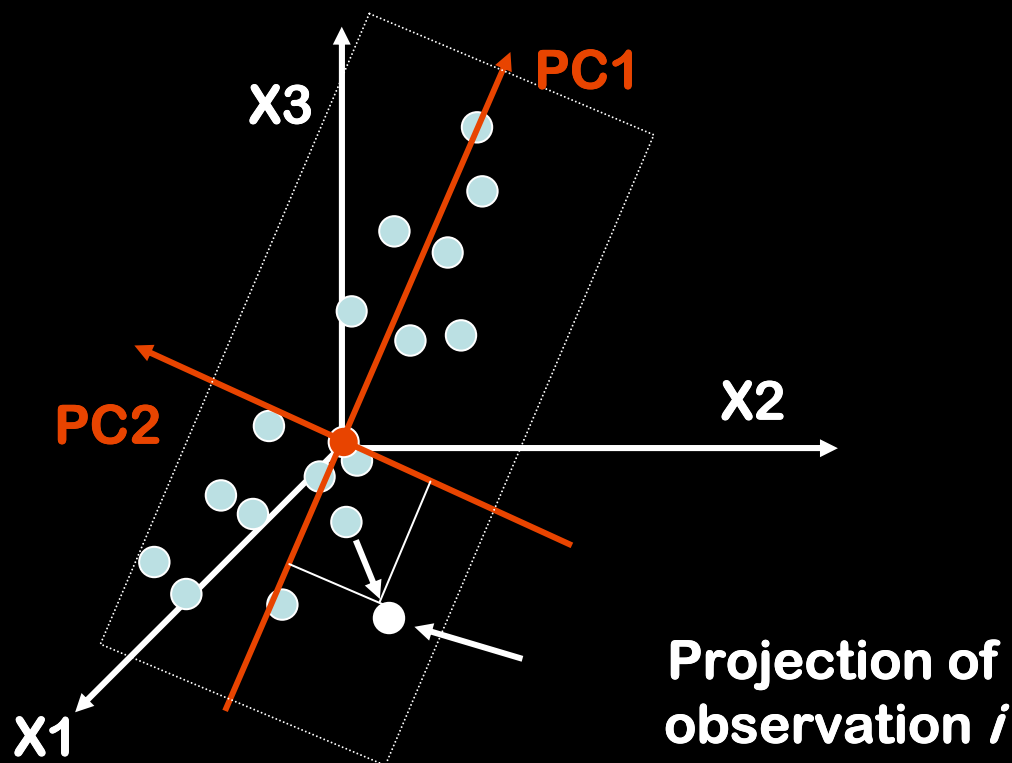
A Geometrical Interpretation of PCA



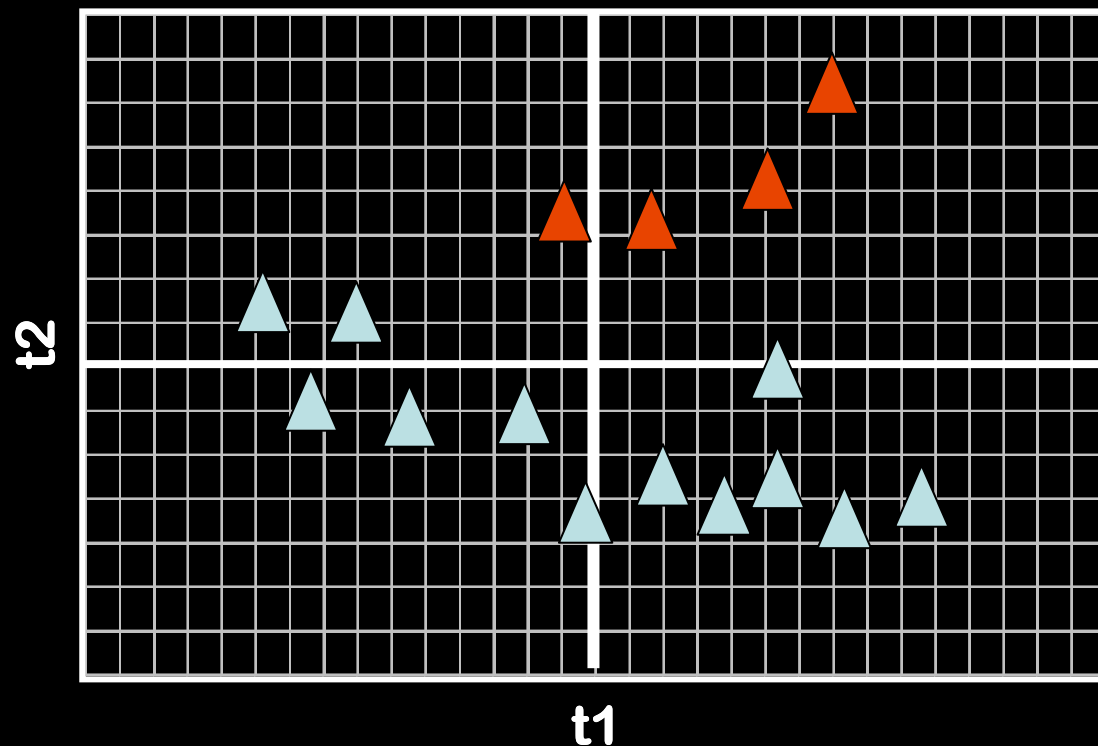
1. Set up k-dimensional space;
2. Plot point;
3. Plot vector of averages at the center of gravity;
4. Mean-center the data;
5. **Generate the first PC:**
Passes through the origin
Best approximates the data in a least squares sense
6. **Generate the second PC:**
Passes through the origin and orthogonal to first PC
Maximally improves the approximation of the X-data

A Geometrical Interpretation of PCA

- First 2 PC's define a plane in the original K-dimensions space.
- By projecting all data points into this plane it is possible to visualize the structure of the data set.
- **Scores (t)** are the coordinates of the original data points in this plane.



Scoring Plot



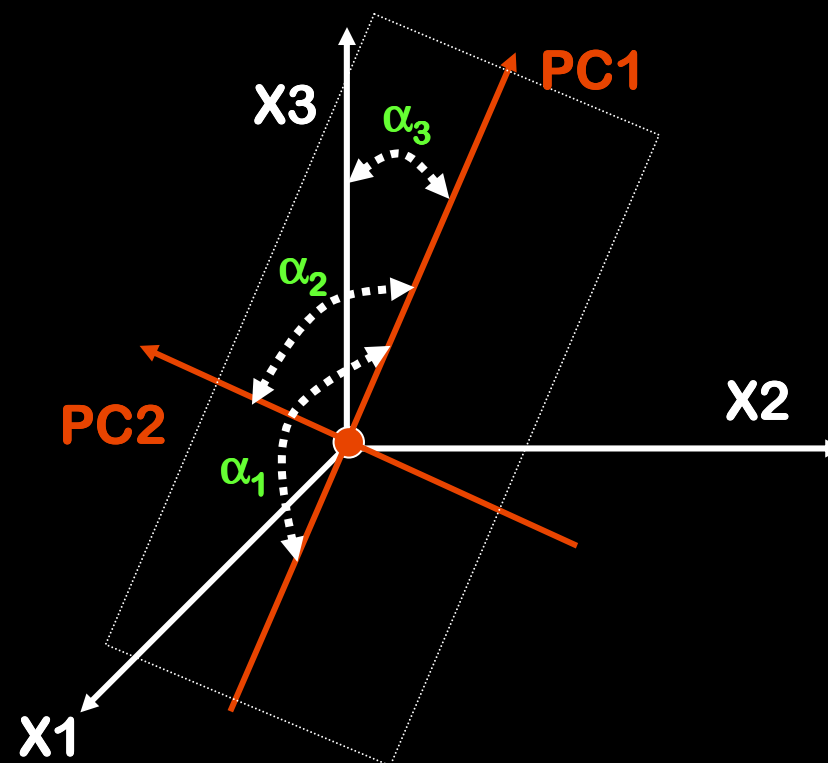
On the score plot,
“Sit together”: similar behavior between descriptors

The Geometrical Meaning of Loading

The loadings express the orientation of the model plane in the original K-dimensional variables space.

The direction of PC1 in relation to the original coordinates is given by the cosine of α_1 , α_2 and α_3 .

With 2 PC's and 3 original variables, 6 loading values (cosine of angles) are needed to specify how the model is positioned in the K-space.



Loading Plot

The contribution (*loading*) of each original variable to each PC.

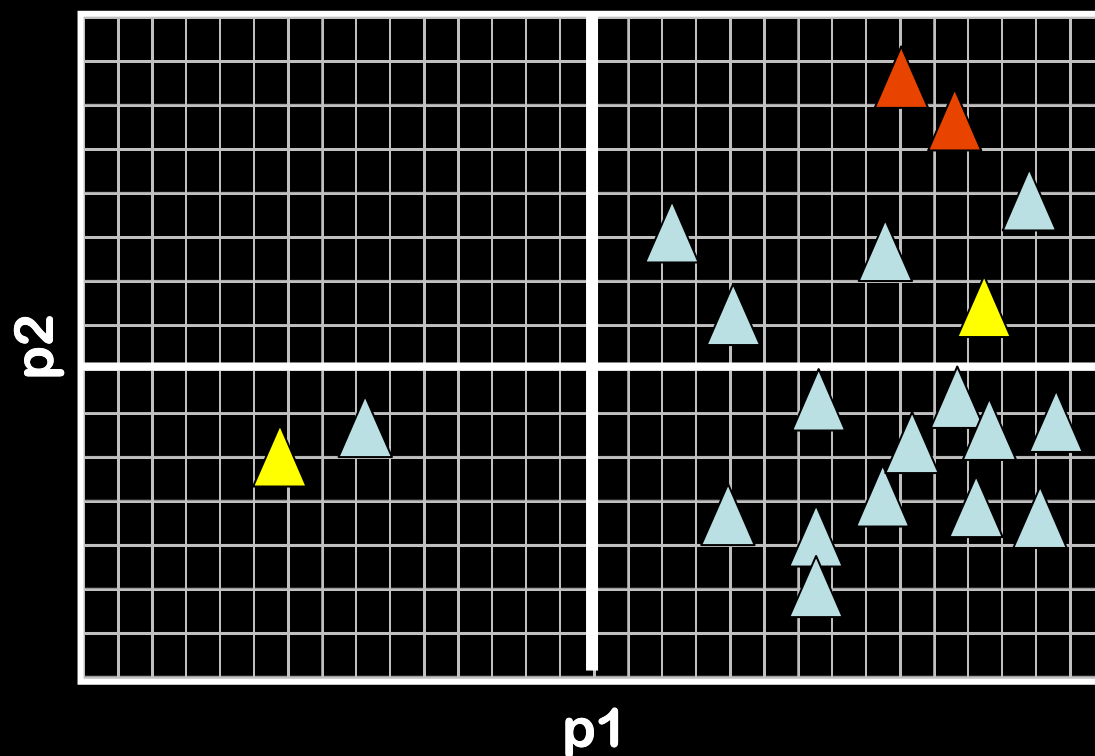
Which variables are responsible for the pattern observed.

PC's can be associated with certain dataset characteristics.

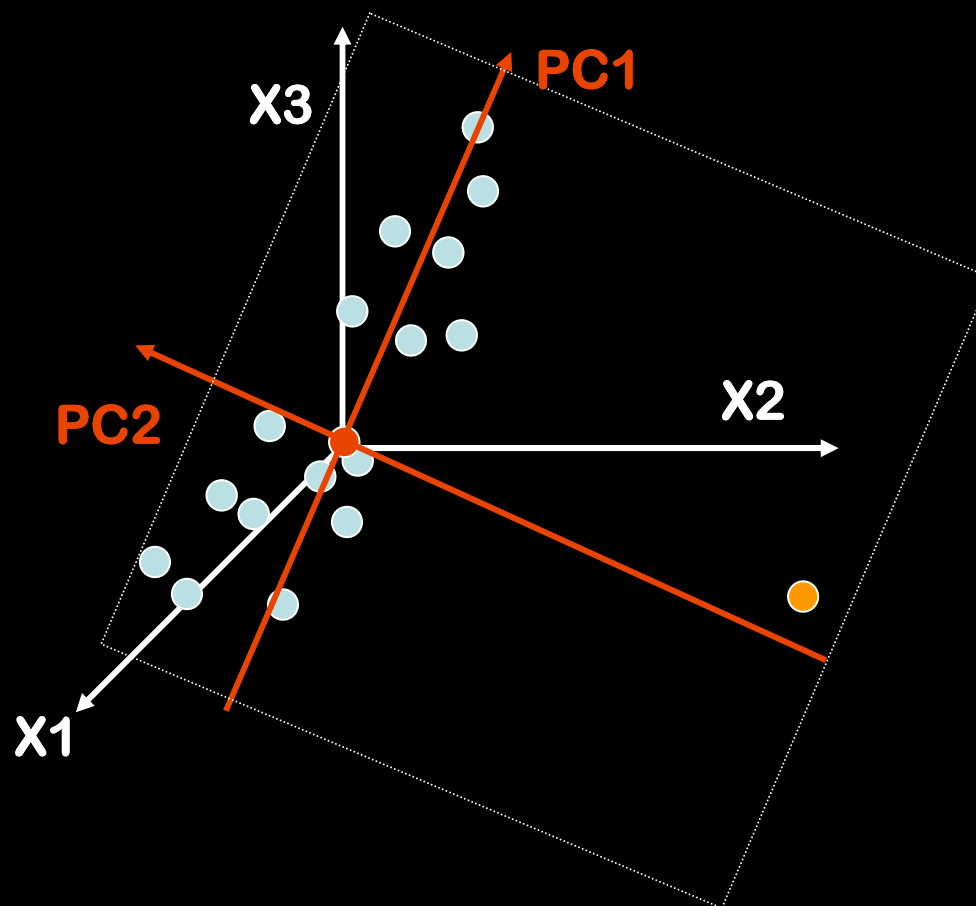
The further away from the origin a variable lies, the stronger impact it has on the model.

Variables correlations:

- ▲ Positively correlated
- ▲ Negatively correlated



Outliers Detection





Back to the real case:

#	LD25	MR	logP	Volume	PM	Surface	density
CH ₂ Cl ₂	0.96	1,62959	1,3044	67,136	84,933	176,1169	1,63677
CHCl ₃	1.45	2,01731	1,7381	75,751	119,378	186,8237	2,03576
CCl ₄	1.53	2,35508	2,4212	83,27	153,823	194,0565	2,3224
CF ₃ CHBrCl	1.31	2,35642	2,3611	88,098	197,381	206,6438	2,85284
CHCl ₂ CHCl ₂	2.42	2,92829	2,4947	94,238	167,85	215,3294	2,26061
Cl ₂ C=CHCl	2.26	2,46705	2,2884	115,83	131,389	241,6985	1,42863
CCl ₂ =CCl ₂	2.26	2,82835	3,3747	132,11	165,834	257,2367	1,46129

$$\log (1/LD_{25}) = a \text{ MR} + b \log P + c \text{ Vol} + d \text{ PM} + e \text{ Sur} + f \text{ dens} + g \text{ X-atom} + h$$

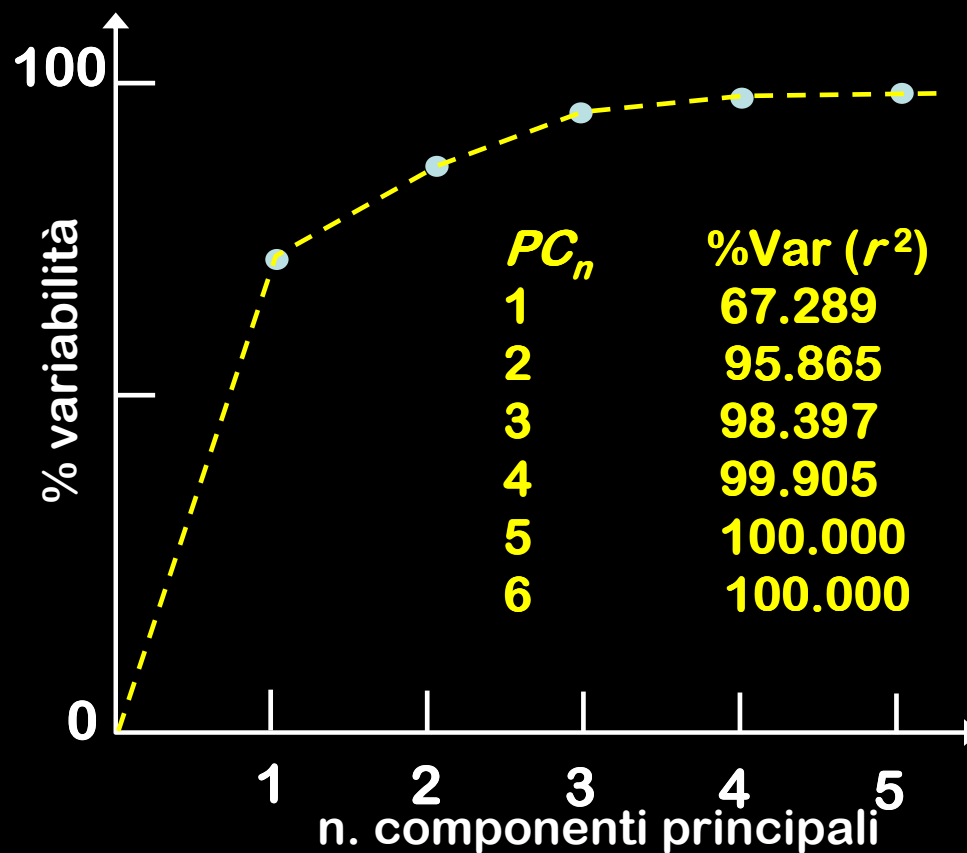
or

$$PC_1 = 0.47 \text{ MR} + 0.34 \log P + 0.008 \text{ Vol} + 0.005 \text{ PM} + 0.006 \text{ Sup} + 0.02 \text{ dens} + 0.15 \text{ atom}$$

$$PC_2 = -0.01 \text{ MR} + 0.08 \log P + 0.01 \text{ Vol} - 0.008 \text{ PM} + 0.008 \text{ Sup} - 1.02 \text{ dens} - 0.18 \text{ g}_2 \text{ atom}$$

$$PC_3 = -4.68 \text{ MR} - 0.67 \log P + 0.03 \text{ Vol} + 0.01 \text{ PM} + 0.02 \text{ Sup} + 0.08 \text{ dens} + 0.62 \text{ g}_2 \text{ atom}$$

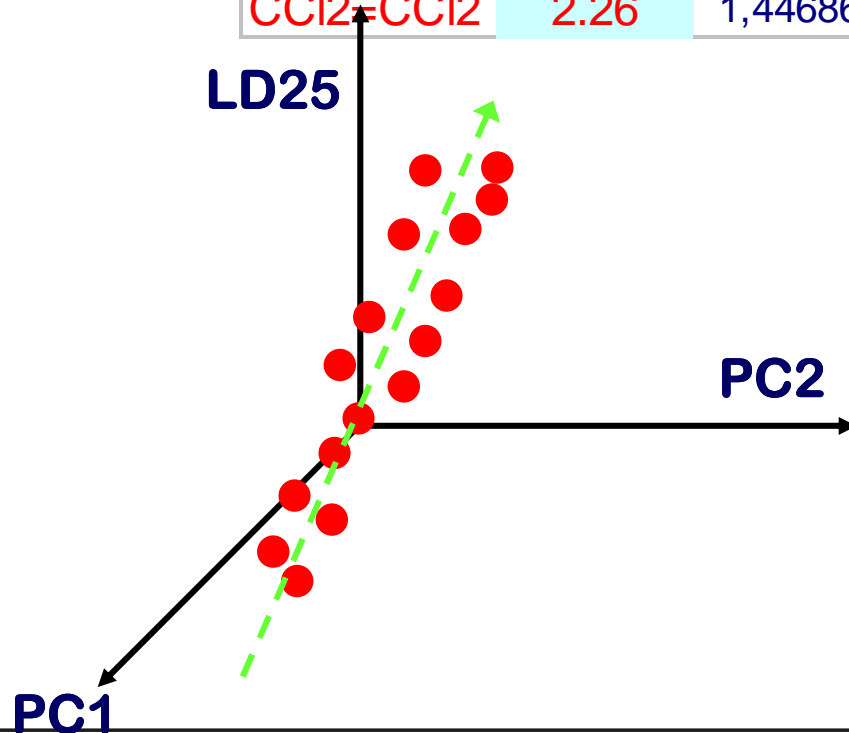
...





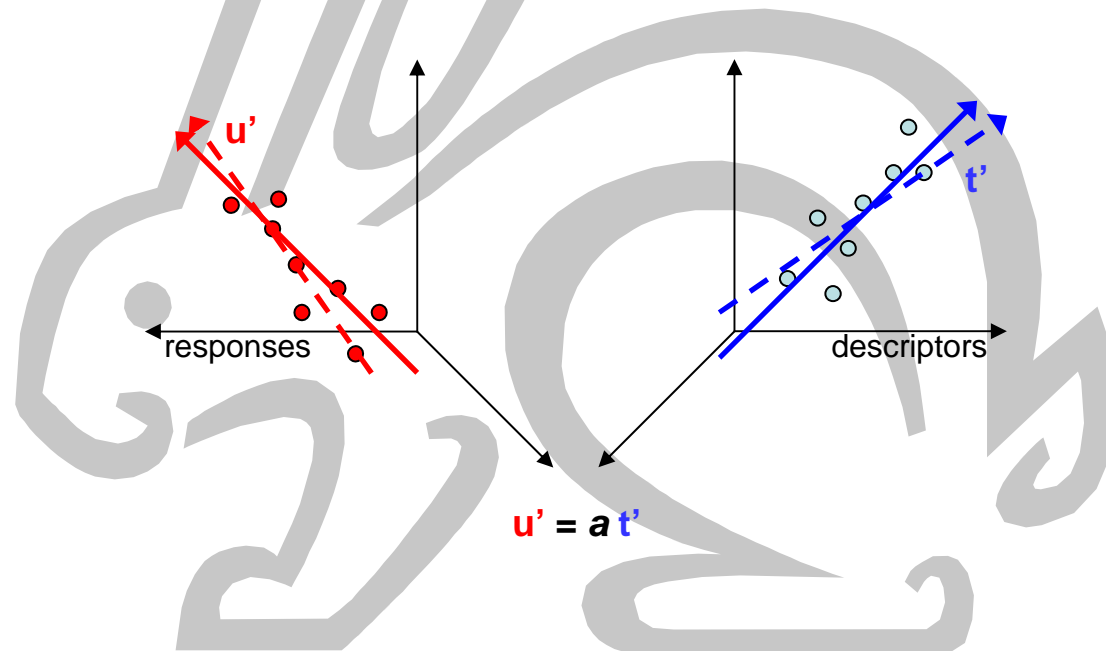
and finally...

#	LD25	PC1	PC2	PC3	PC4	PC5
CH ₂ Cl ₂	0.96	-1,78677	0,583913	0,285292	0,101483	1,469406
CHCl ₃	1.45	-0,97519	-0,07003	-0,20004	0,095932	-1,21311
CCl ₄	1.53	-0,13863	-0,63091	-0,75064	1,407425	-1,00089
CF ₃ CHBrCl	1.31	0,500215	-1,74077	1,586333	-0,26991	0,353284
CHCl ₂ CHCl ₂	2.42	0,616411	-0,55574	-1,81367	-1,19166	0,764962
Cl ₂ C=CHCl	2.26	0,337096	1,253576	0,681746	-1,39385	-1,0778
CCl ₂ =CCl ₂	2.26	1,446869	1,159959	0,210977	1,250575	0,704142

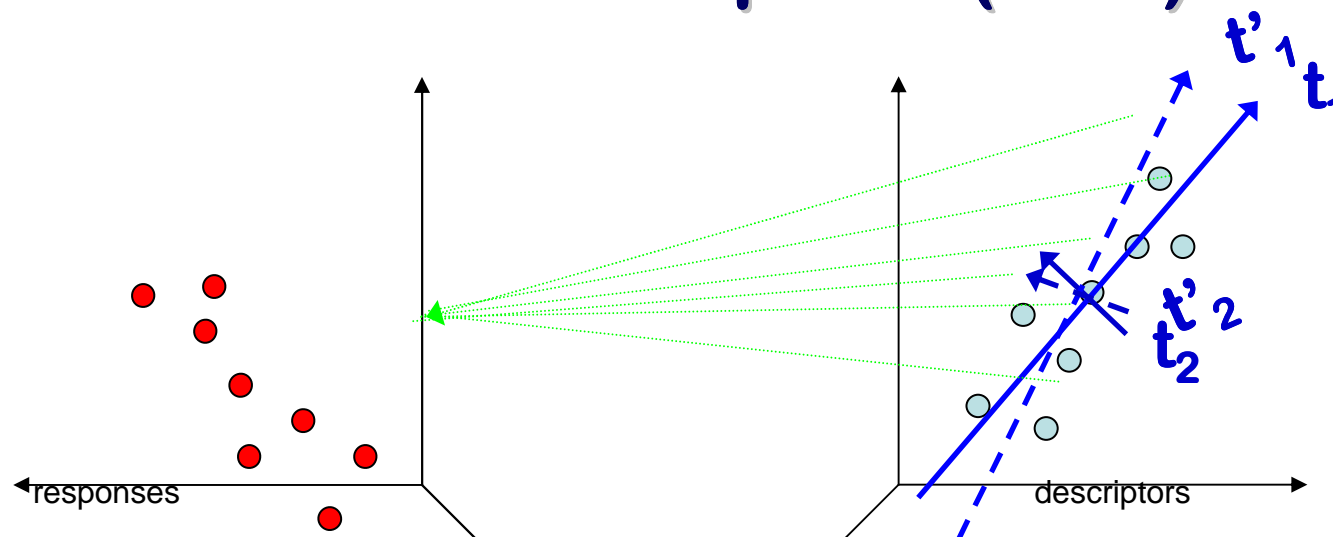


$$\log (1/LD_{25}) = m PC1 + n PC2$$

Partial Least Square Regression



Partial Least Square (PLS)



t'_n : latent variables

$$t'_1 = a_1x_1 + b_1x_2 + \dots + m_1x_m$$

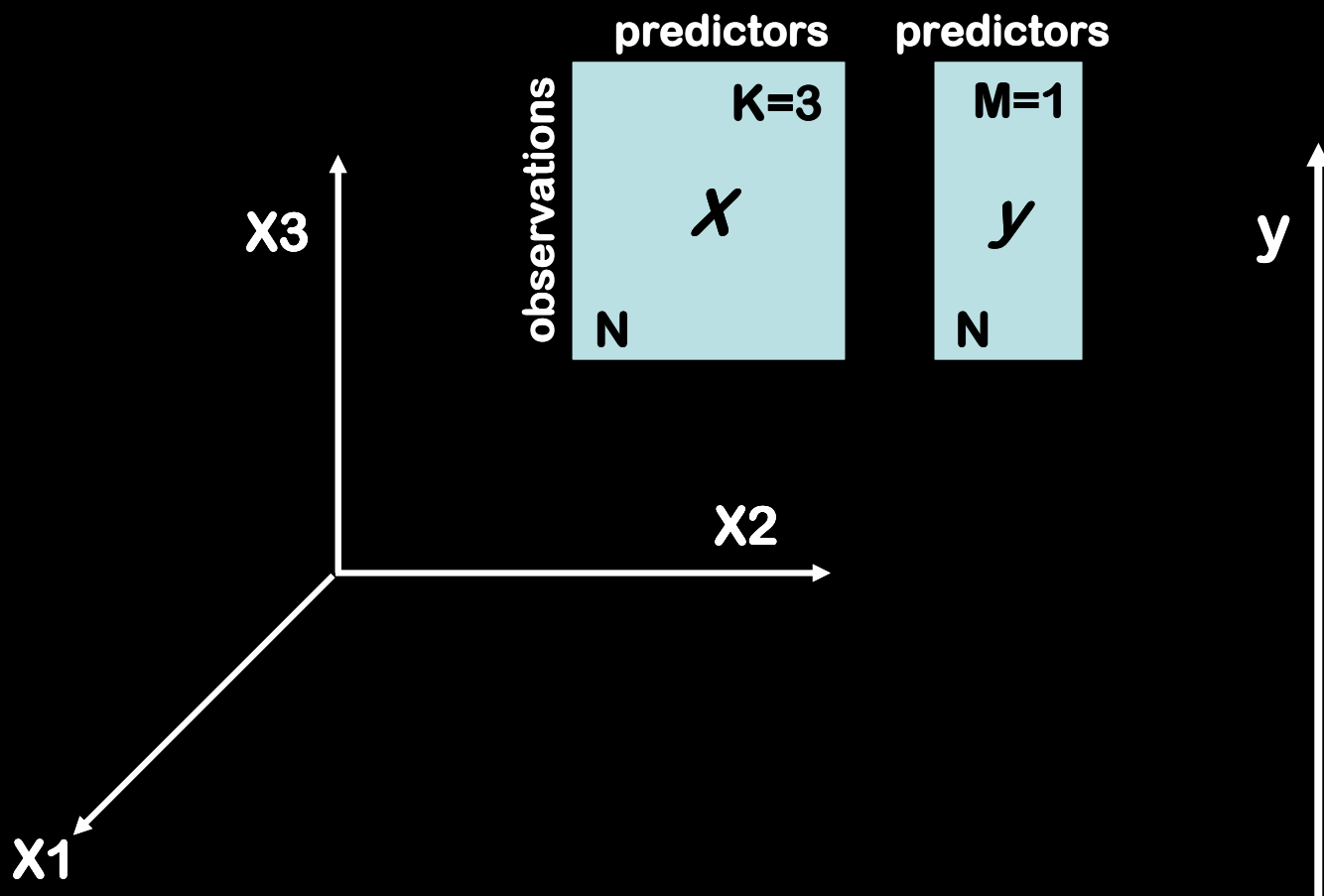
$$t'_2 = a_2x_1 + b_2x_2 + \dots + m_2x_m$$

$$y_1 = a_1t'_1 + b_1t'_2 + \dots + m_1t'_m$$

...

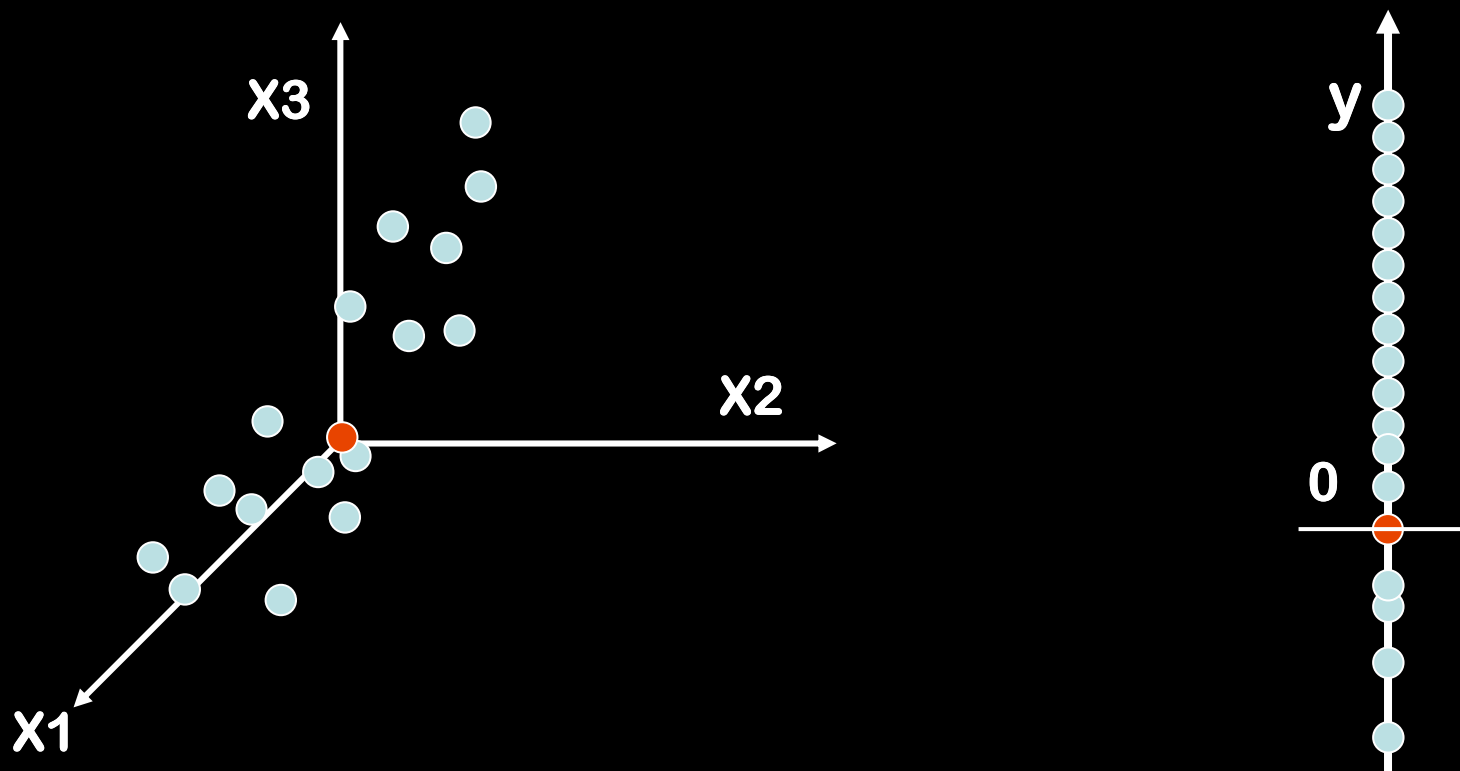
$$y_n = a_nt'_n + b_nt'_n + \dots + m_nt'_m$$

A geometrical Interpretation of PLS



A geometrical Interpretation of PLS

PLS describes the relationship between the (mean centered) positions of the observations in the predictor space (x) and their (mean centered) positions in the response space (y).



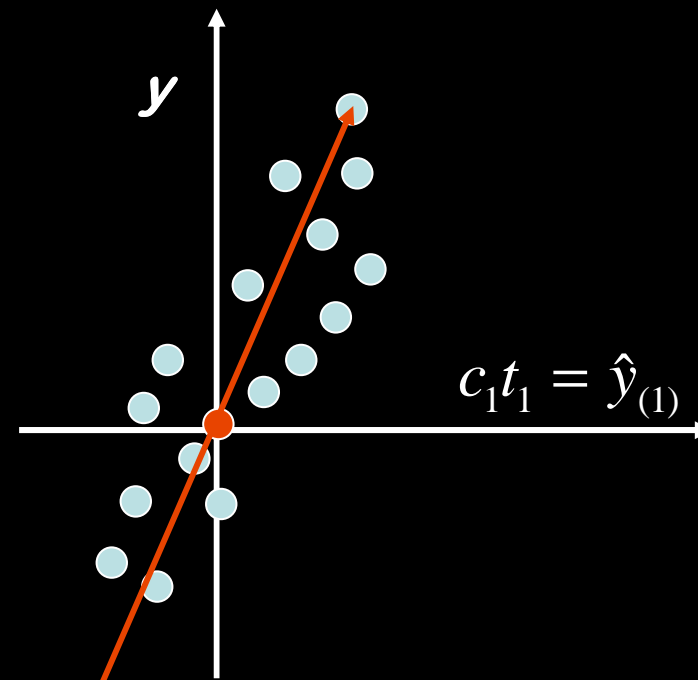
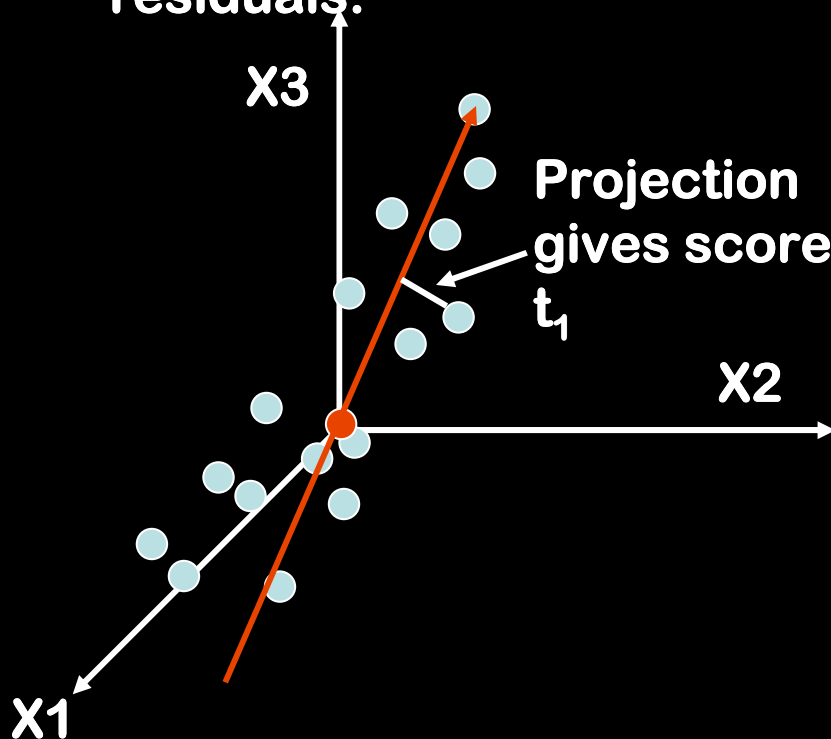
The First PLS Component

Remember our goal:

- Predict y from X .
- Since there are too many X 's, generate a new set of variables in a PCA-manner.
- Predict y from a new variable t , t being a linear combination of the original X 's.
- PLS assumes y predicted from t in a linear fashion, $y = at + b$.

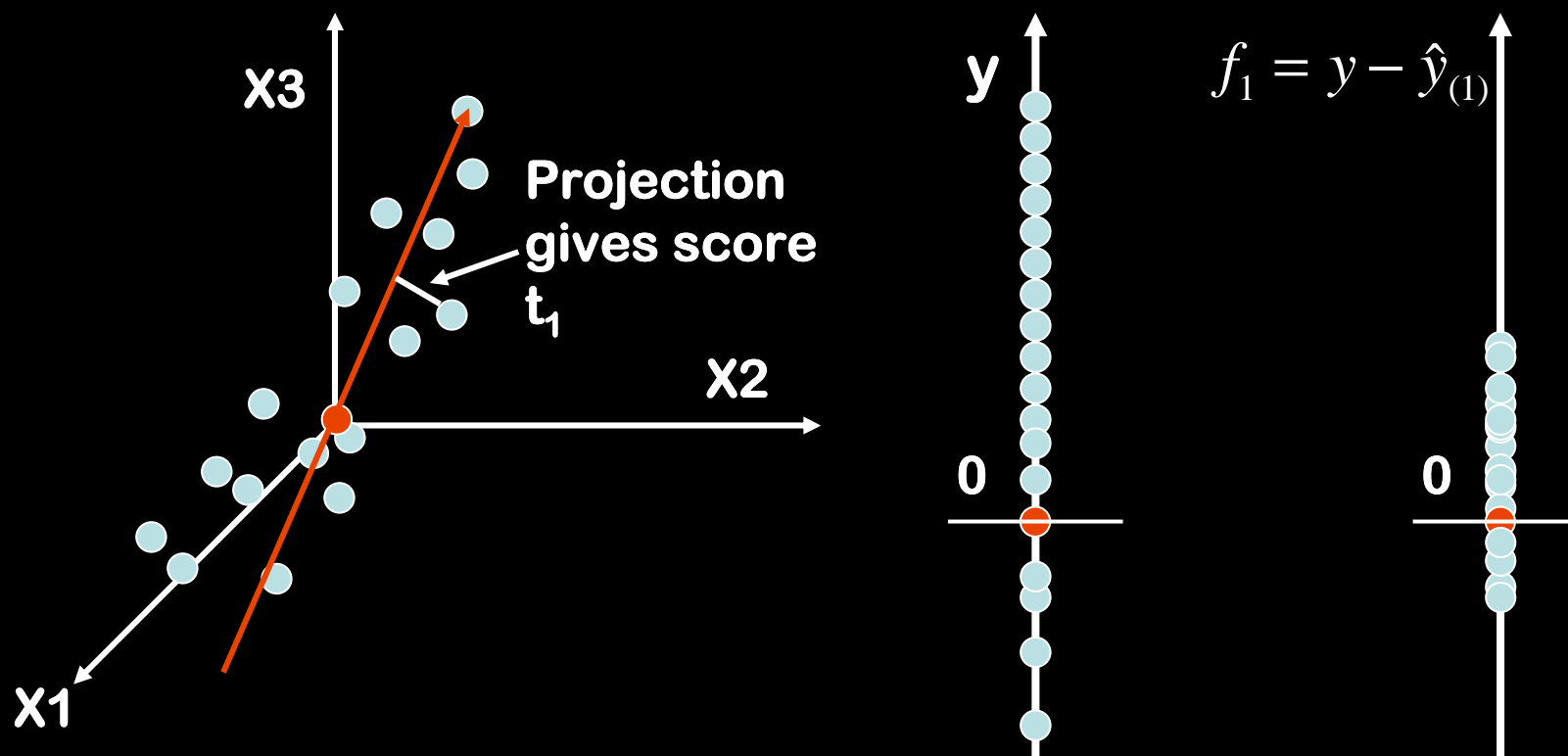
The First PLS Component

- Line in the X-space which well approximates the points in a least squares manner and at the same time provides a good correlation with the y vector.
- Degree of correlation with the y vector is judges from the residuals.



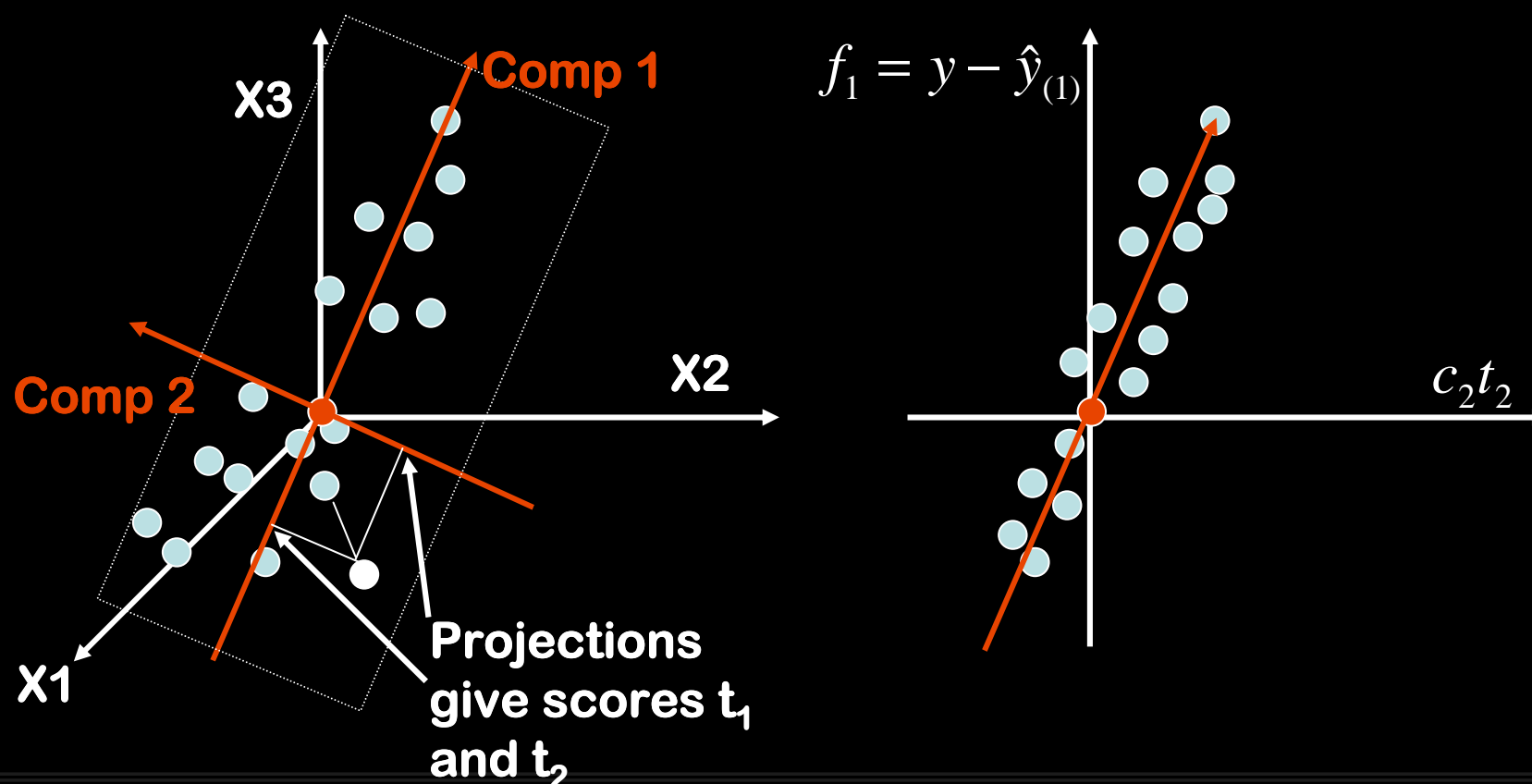
Residual of the First PLS Component

- The residual vector is shorter than the original y vector indicating that the first PLS component (first latent variable) accounted for a large part of the variation in y .



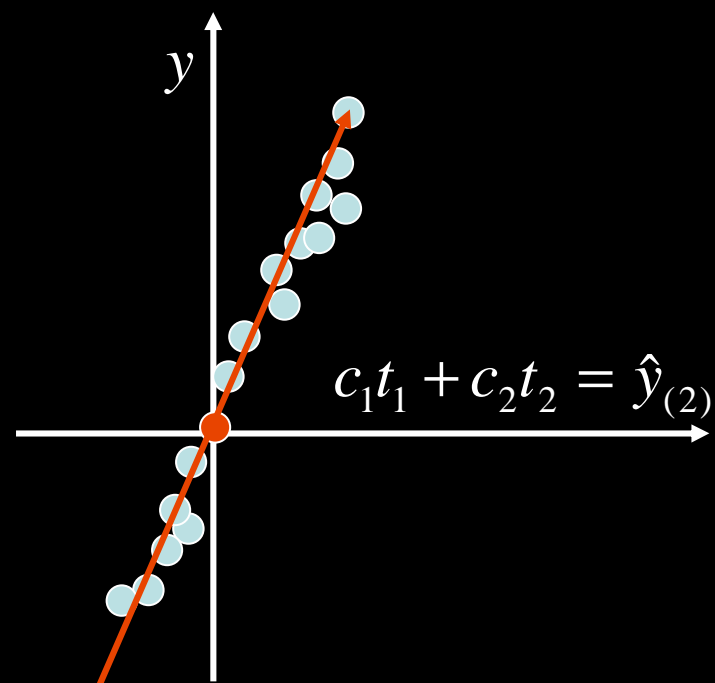
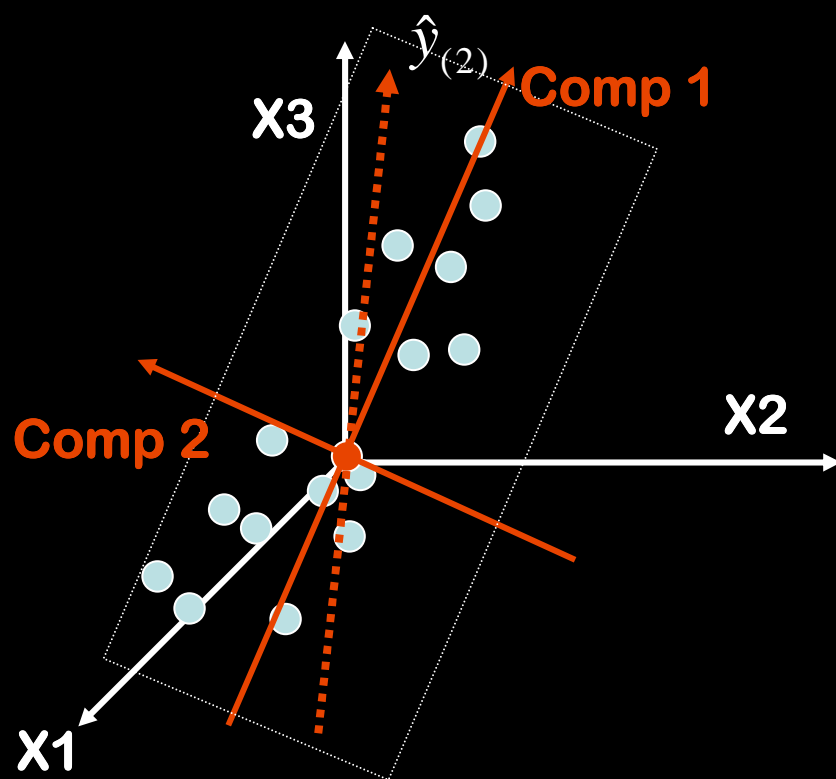
The Second PLS Component

- Line in the X-space which improves the description of the X data and at the same time provides a good correlation with the f_1 vector.
- Orthogonal to the first PLS component



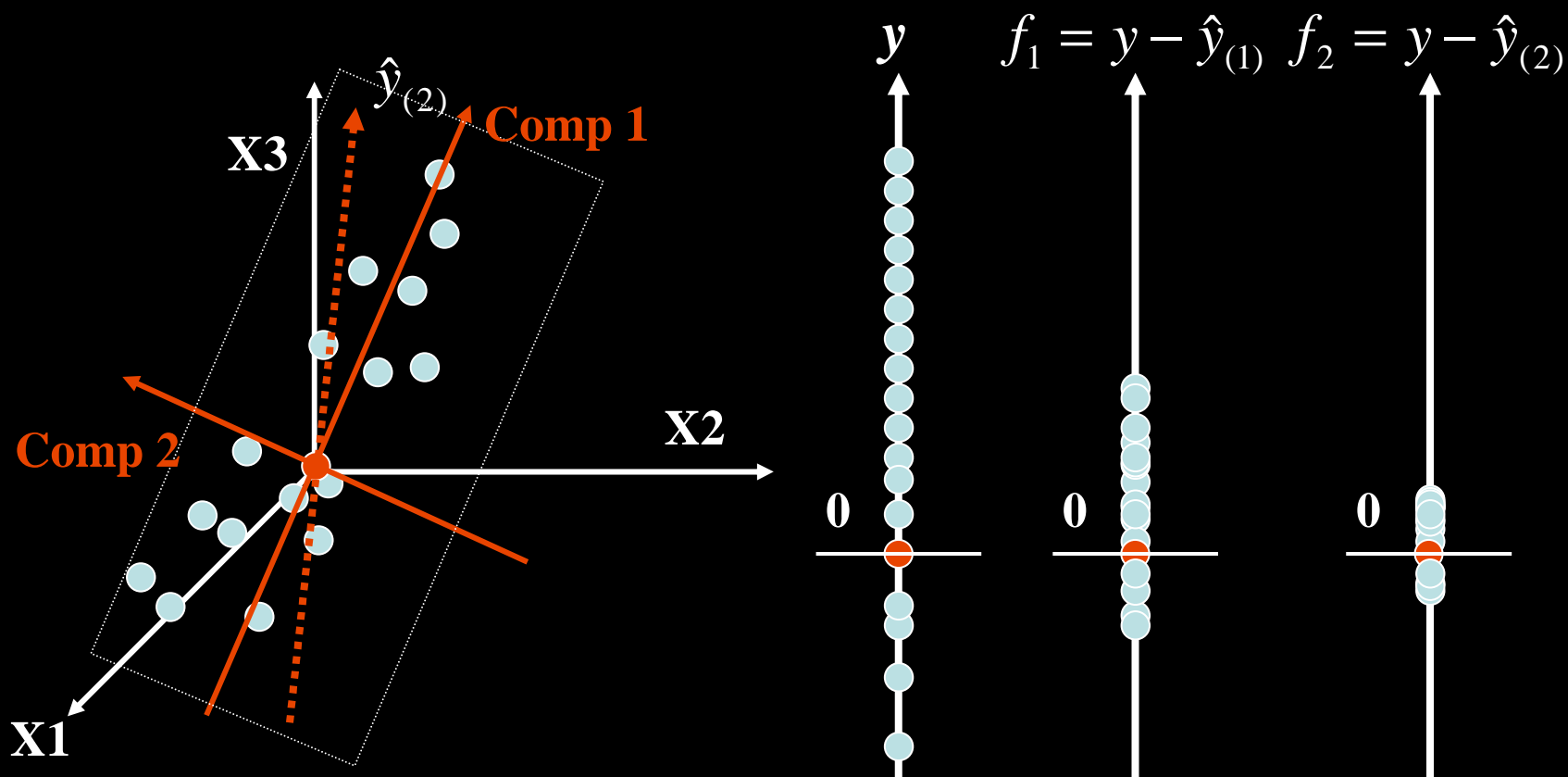
The Cumulative Effect of two PLS Components

- $y_{(2)}$ is a vector addition of the two first components in the X space.
- $y_{(2)}$ is a better predictor of the y values.

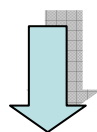
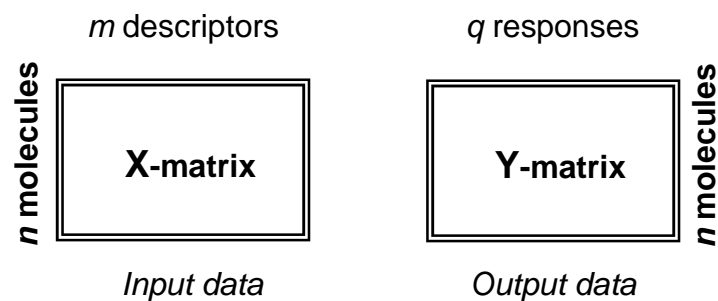


Explanatory Power of PLS

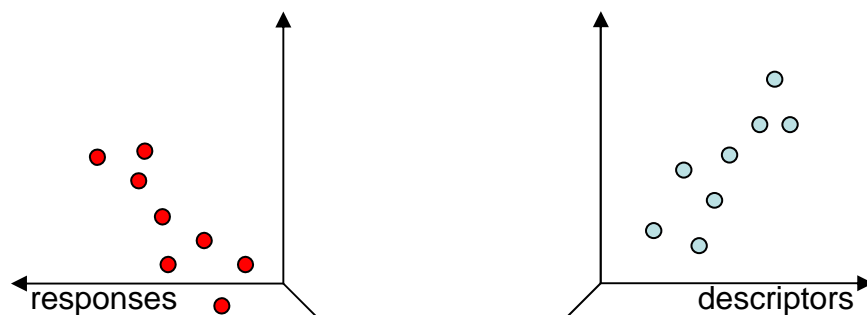
- $y_{(2)}$ is a vector addition of the two first components in the X space.
- $y_{(2)}$ is a better predictor of the y values.



Linear Methods:



Multiple Regression Analysis (MRS)

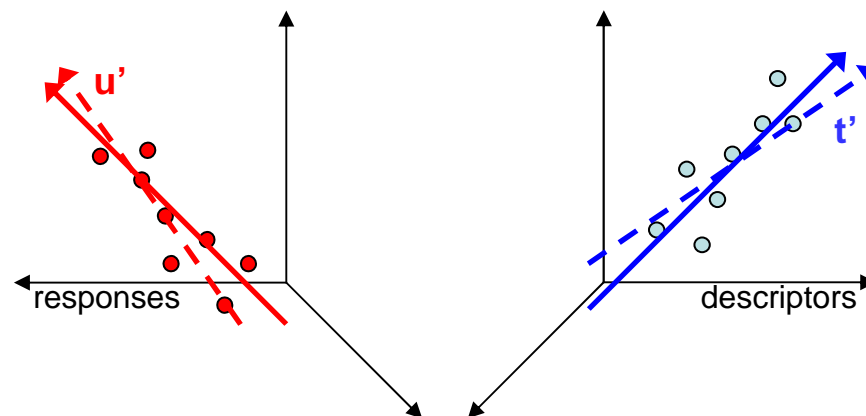


$$y_1 = a_1 x_1 + b_1 x_2 + \dots + m_1 x_m$$

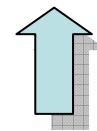
...

$$y_n = a_n x_1 + b_n x_2 + \dots + m_n x_m$$

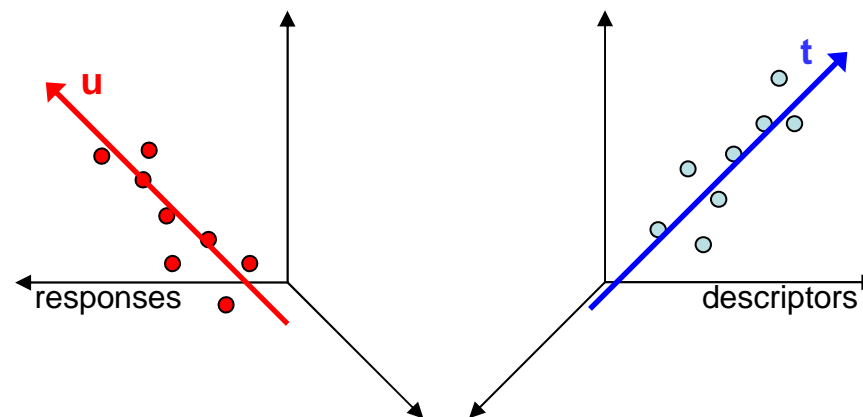
Partial Least Square (PLS)



$$u' = a t'$$



Principle Components Regression (PCR)



$$u = a t$$

Artificial Neural Networks (ANN): una soluzione alternativa alla forzata linearità delle QSARs.

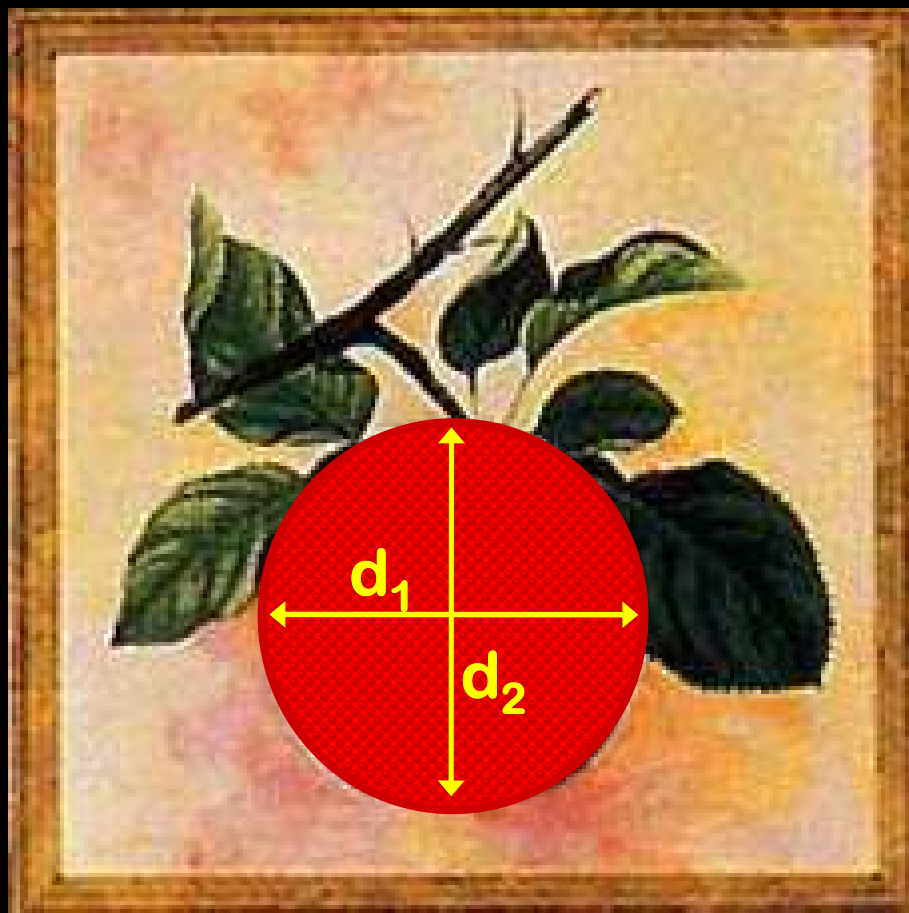
Artificial Neural Networks (ANN)



How we can distinguish these two objects?

Experience → Multiple Experience → Category

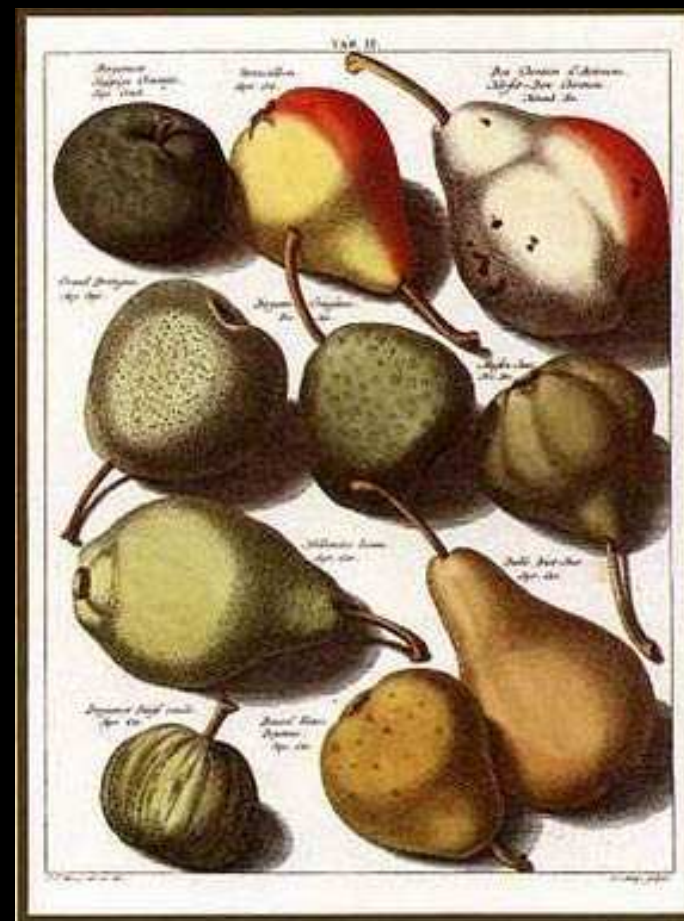
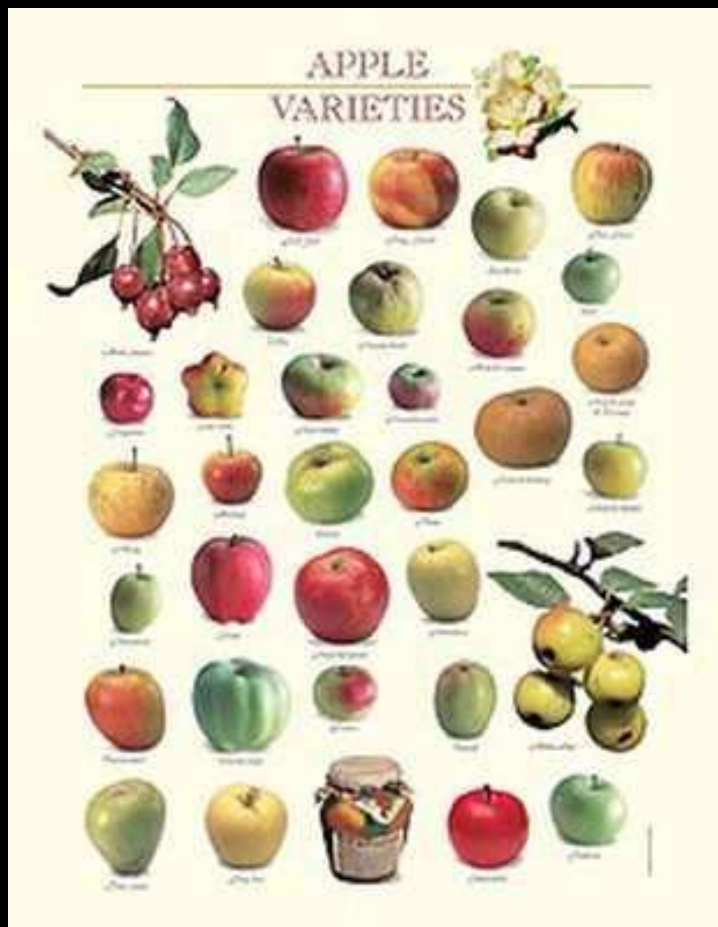
Artificial Neural Networks (ANN)



Color
Symmetry
...

Why apple?

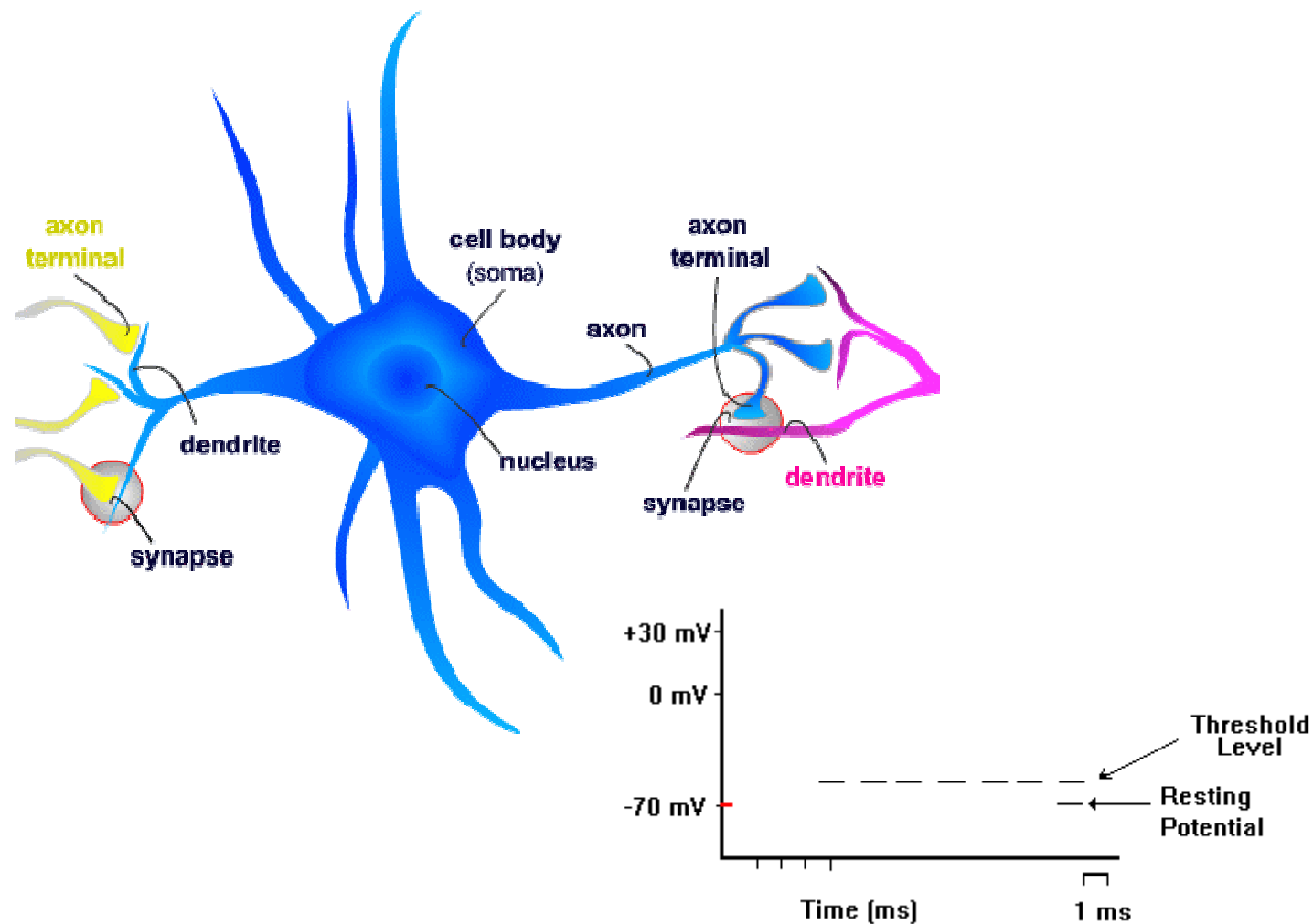
Artificial Neural Networks (ANN)



How we can distinguish apples by pears?

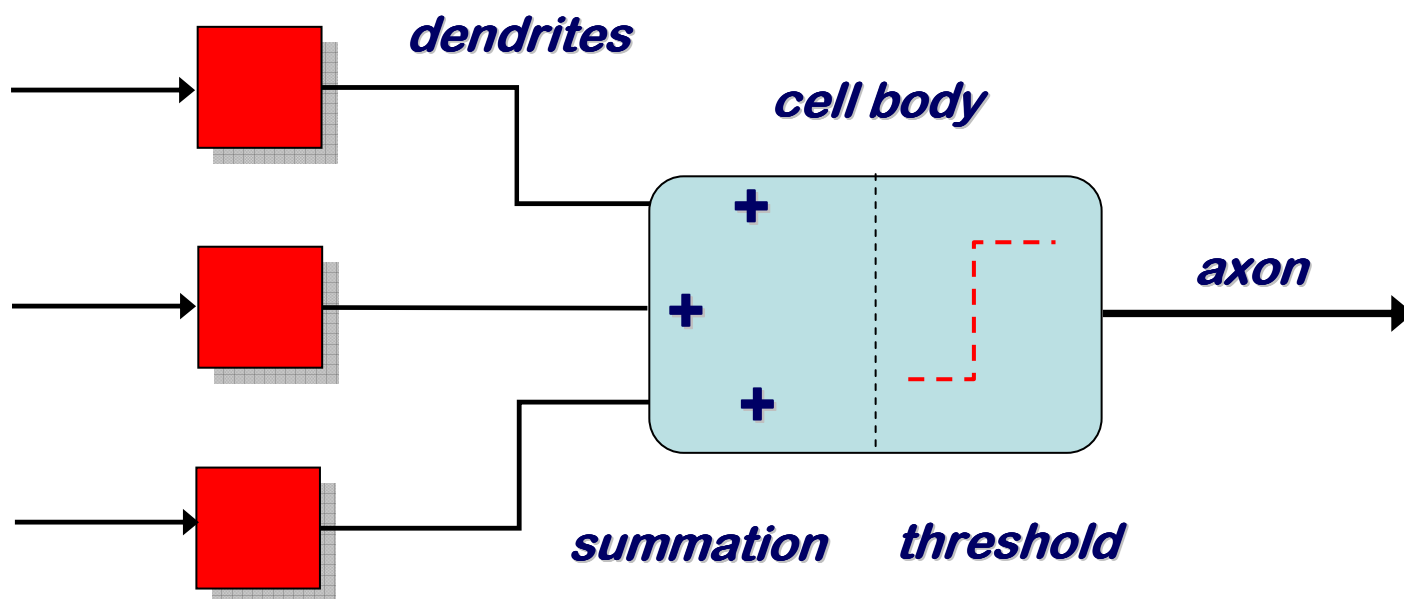


Do you remember the structure of neurons?





From human neurons to artificial neurons...



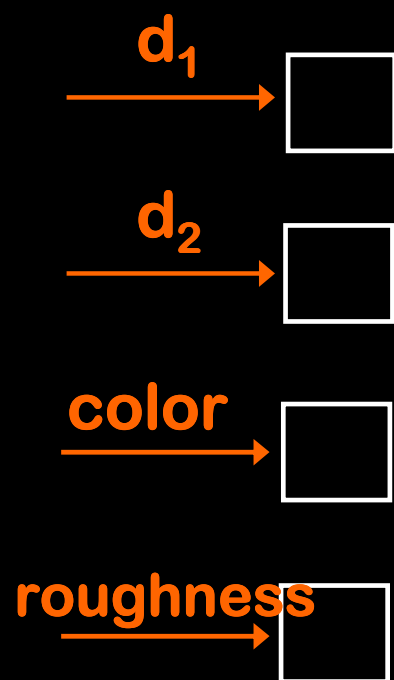
Artificial Neural Networks (ANN):

Simplify structure of an ANN

INPUT

NEURAL NETWORKS

OUTPUT

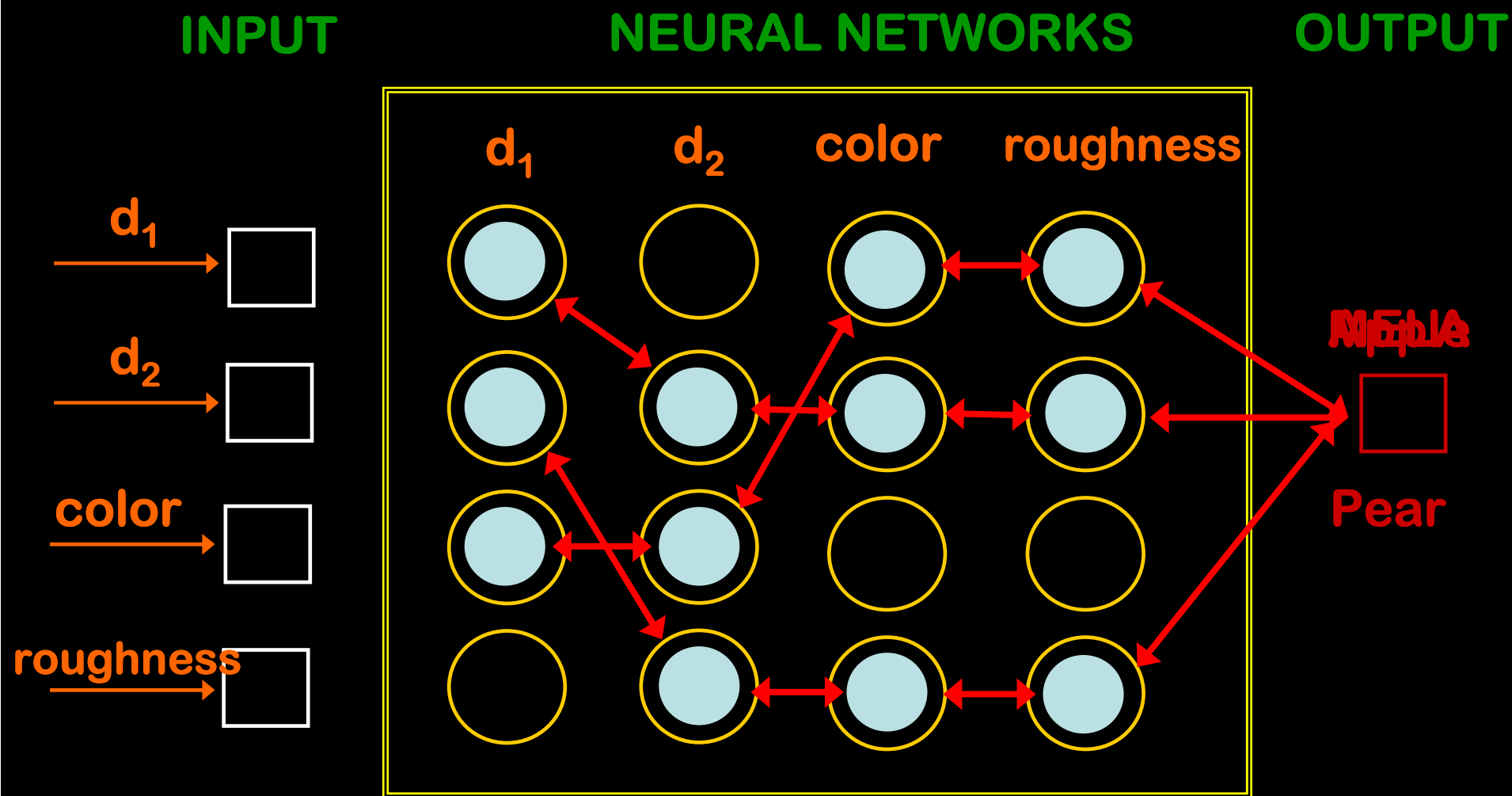


d_1	d_2	color	roughness
5-7	5-9	yellow	--
7-9	9-12	red	-
9-11	12-16	green	+
11-13	16-20	brown	++



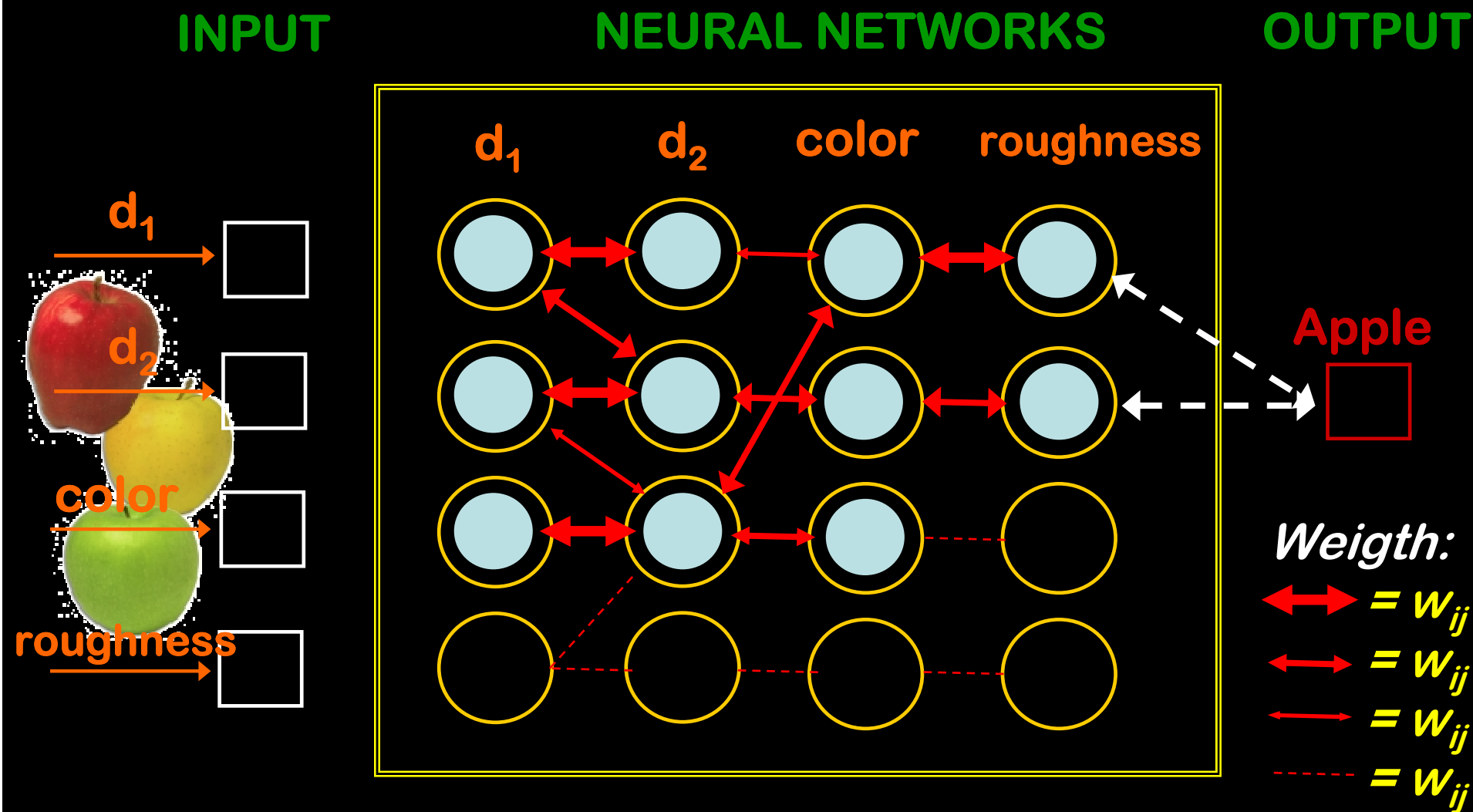
Artificial Neural Networks (ANN):

Phase 1 – Learning.



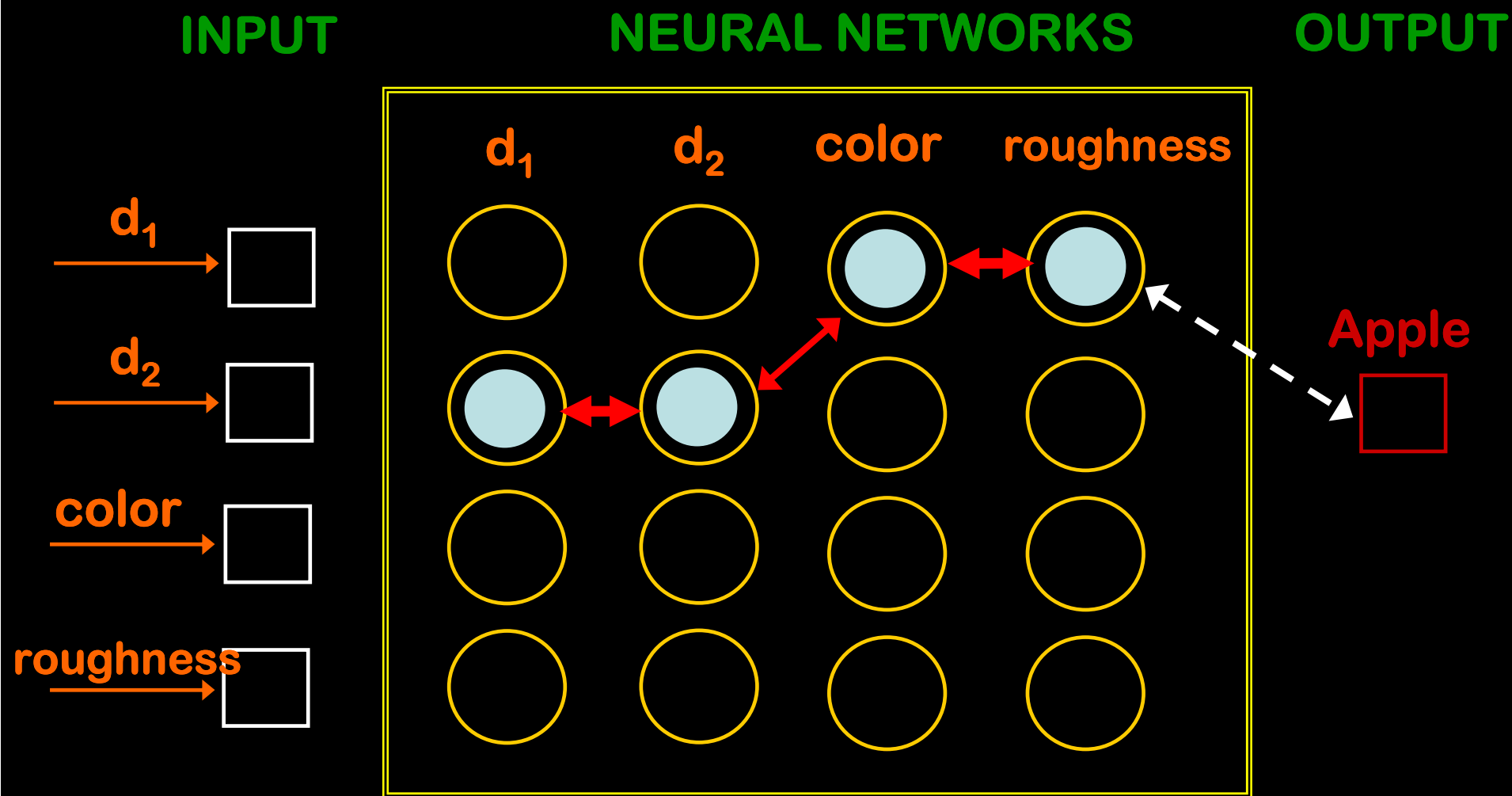
Artificial Neural Networks (ANN):

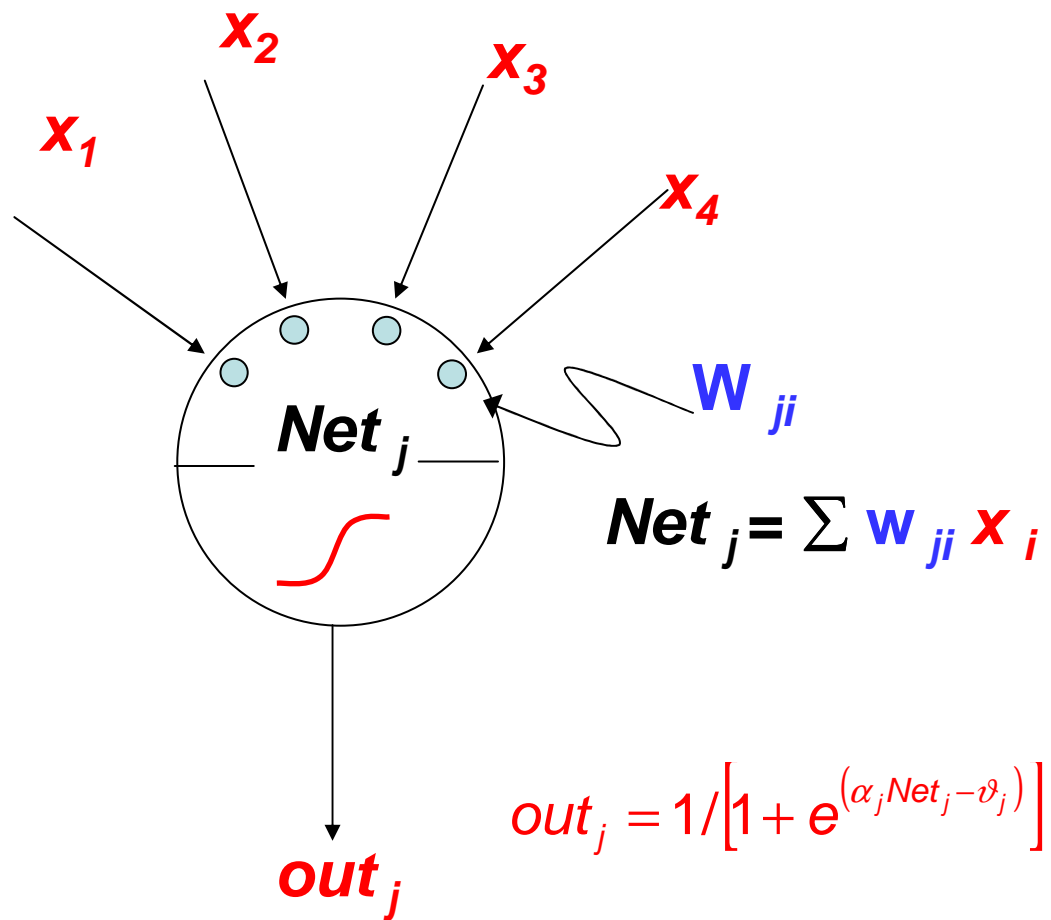
Phase 1 – Learning.



Artificial Neural Networks (ANN):

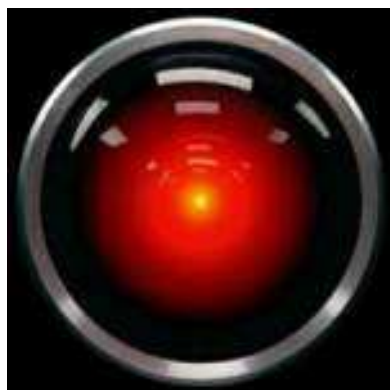
Phase 2 – Recognition.







A view of HAL 9000's Main Terminal



HAL's iconic camera eye.

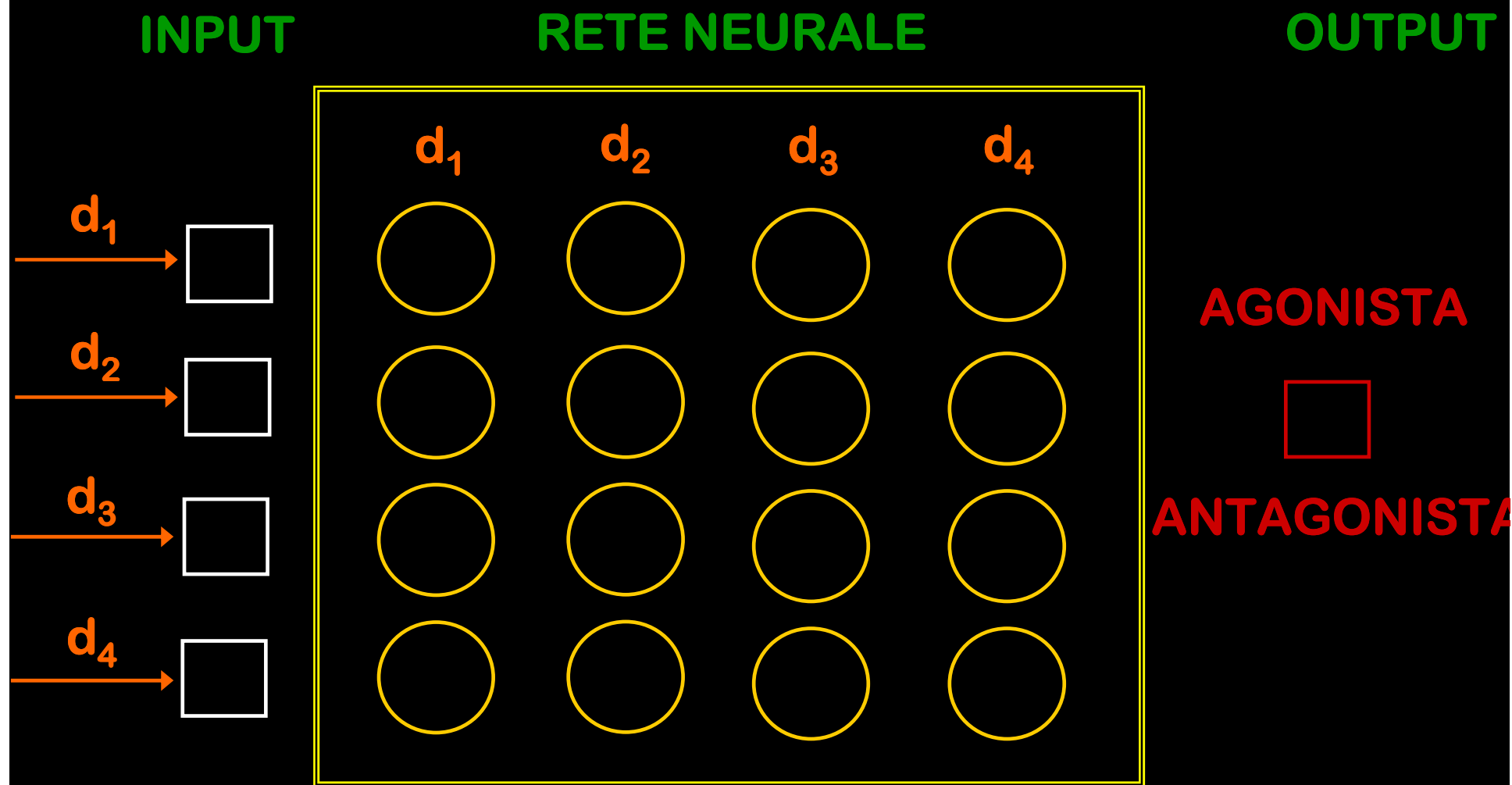


A view of HAL 9000's Central Core in the *Discovery*.

Let me put it this way, Mr. Amer. The 9000 series is the most reliable computer ever made. No 9000 computer has ever made a mistake or distorted information. We are all, by any practical definition of the words, foolproof and incapable of error.

Artificial Neural Networks (ANN):

Applicazione in chimica farmaceutica.





Now, it is time to back at the nature of molecular descriptors.

