



спецкурс
“СУЧАСНІ КОМП’ЮТЕРНІ МЕТОДИ В ХІМІЇ”

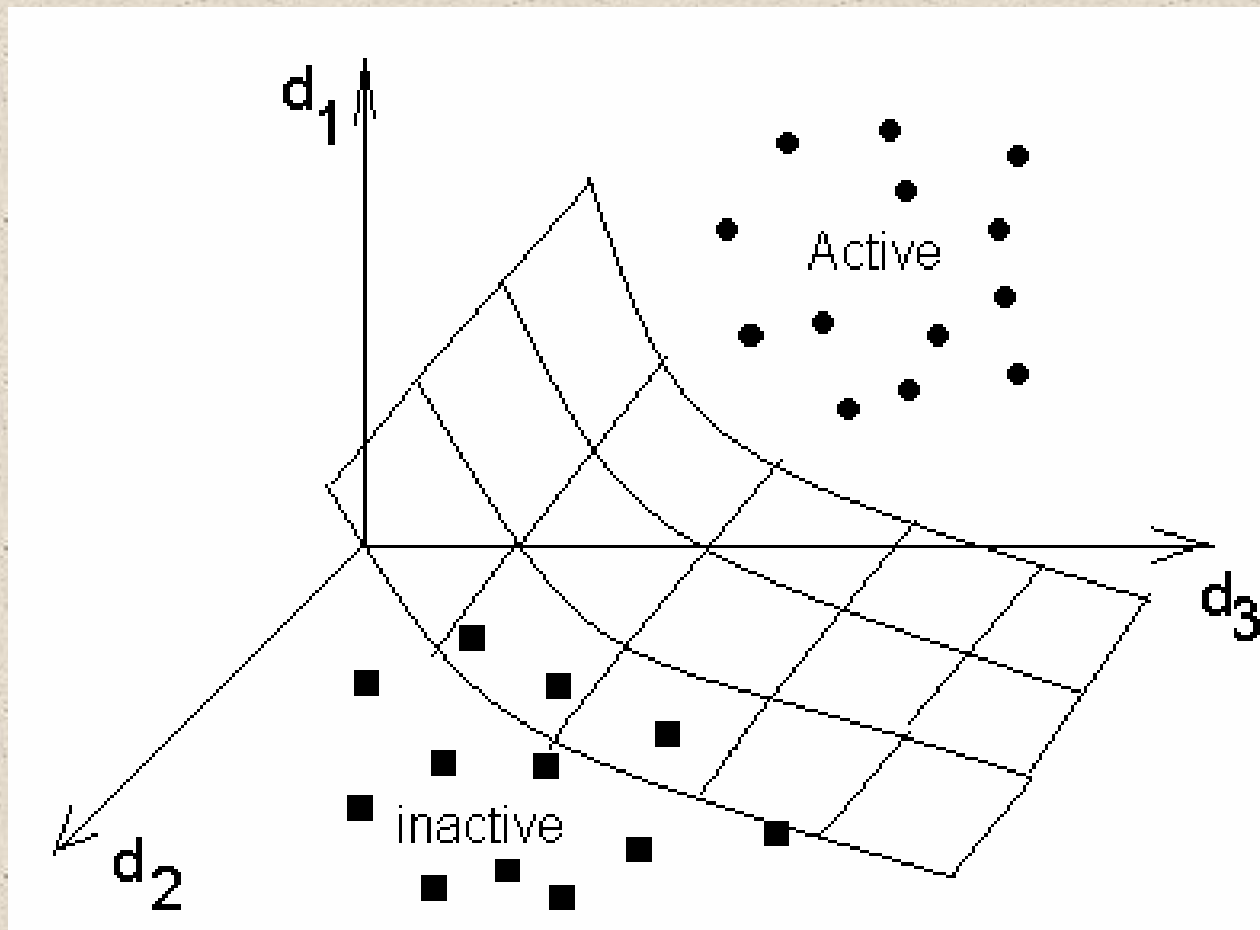
Методи класифікації

В. В. Іванов

Materials Chemistry Department
V. N. Karazin National University,
61077, Kharkiv, Ukraine
ivv34@yahoo.com

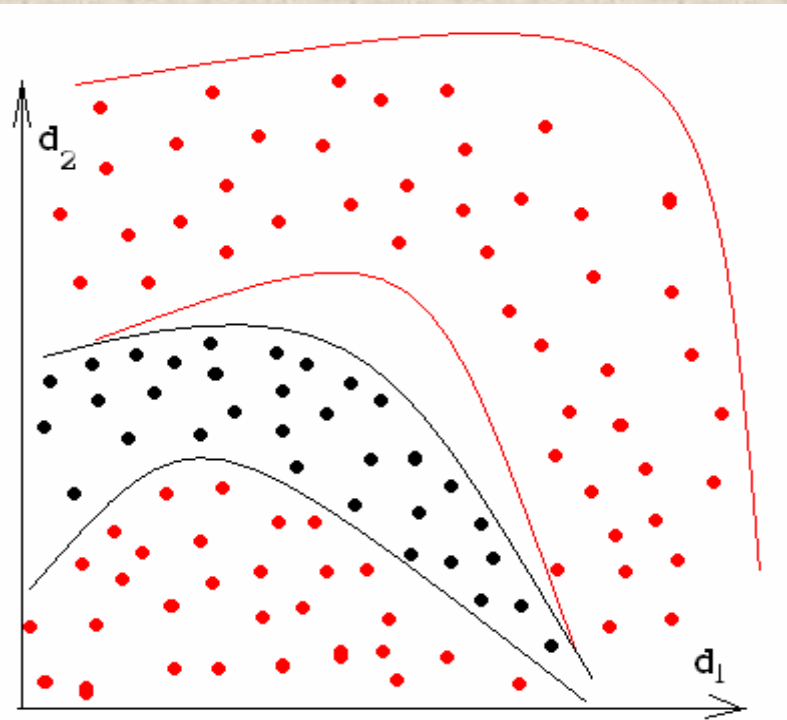
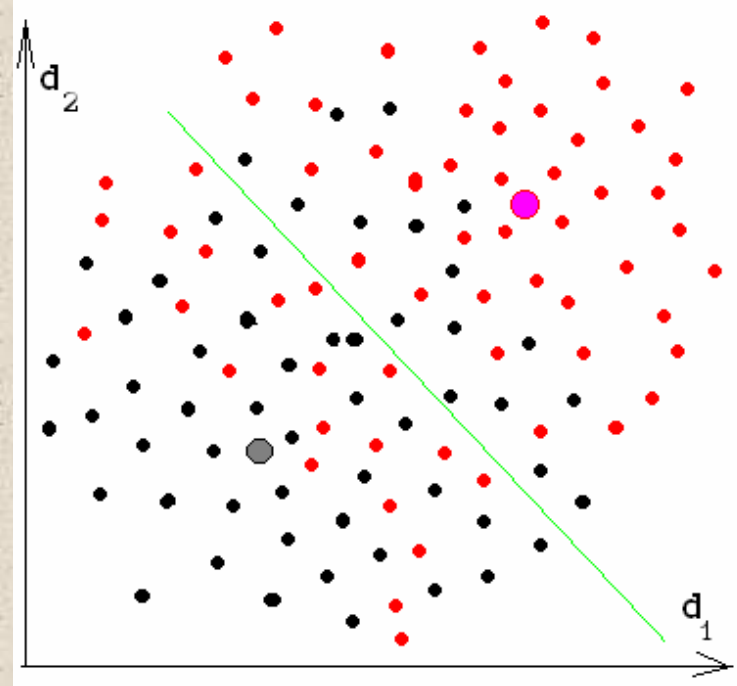
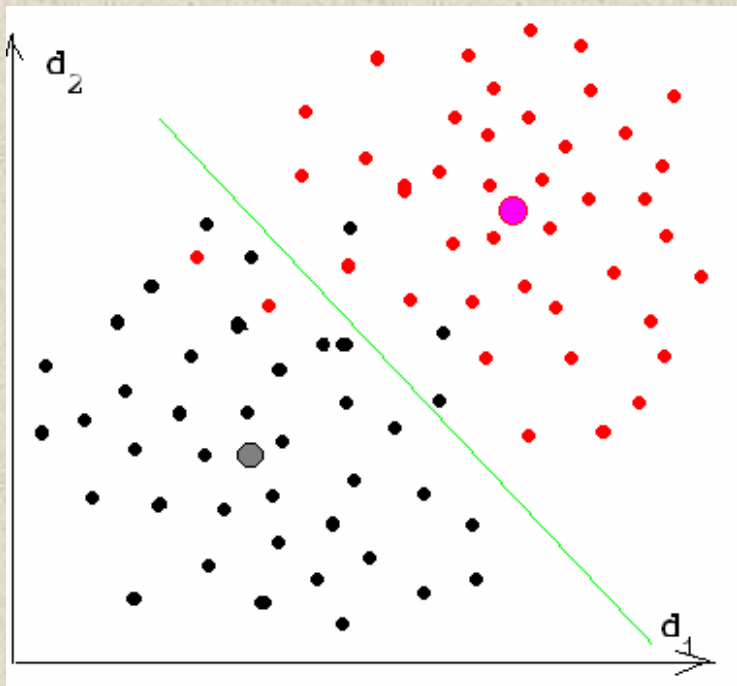
ДИСКРИМІНАНТНИЙ (дискримінаційний) АНАЛІЗ

$$D = a_0 + a_1 d_1 + a_2 d_2 + a_3 d_3 + b_{11} d_1^2 + b_{12} d_1 d_2 + b_{13} d_1 d_3 \dots$$



Мета класифікації

- Ідентифікація змінних (дескрипторів), які Найкращім чином розділяють *активні* і *неактивні* молекули
- Створення *класифікаційного правила* для розділення систем на групи



Метод Фішера лінійна дискримінація (LDA) по 2 класам

$$\mathbf{x} = (d_1, d_2, \dots, d_p)$$

$$D = a_1 \cdot d_1 + a_2 \cdot d_2 + \dots$$

$$W_{ij} = \sum_{\alpha=1}^2 \sum_{m=1}^{n_{\alpha}} (d_{i\alpha m} - \bar{d}_{i\alpha})(d_{j\alpha m} - \bar{d}_{j\alpha})$$

$$C_{ij} = \sum_{\alpha=1}^2 \sum_{m=1}^{n_{\alpha}} (d_{i\alpha m} - \bar{d}_i)(d_{j\alpha m} - \bar{d}_j)$$

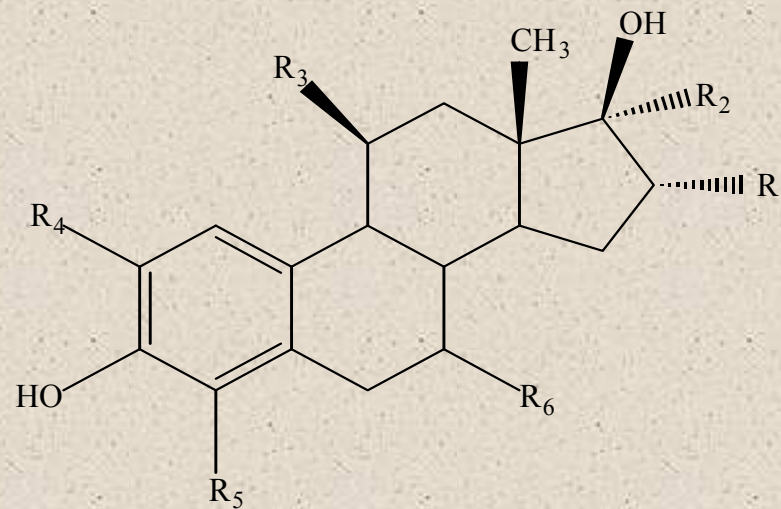
$$T_{ij} = C_{ij} - W_{ij}$$

$$\max \lambda(\vec{a}) = \frac{\sum_{i,j} a_i T_{ij} a_j}{\sum_{i,j} a_i W_{ij} a_j} \quad 5$$



1938, Рональд Фішер

Дискримінантний аналіз Активність естрадіолів



Активність:
 $\text{Log(RBA)} > 1.5$

$$1.032 - 1.49 \varepsilon_{\text{HВМО}} - 0.089 \alpha + 1.41 \lg P - 0.17(\lg P)^2 > 0$$

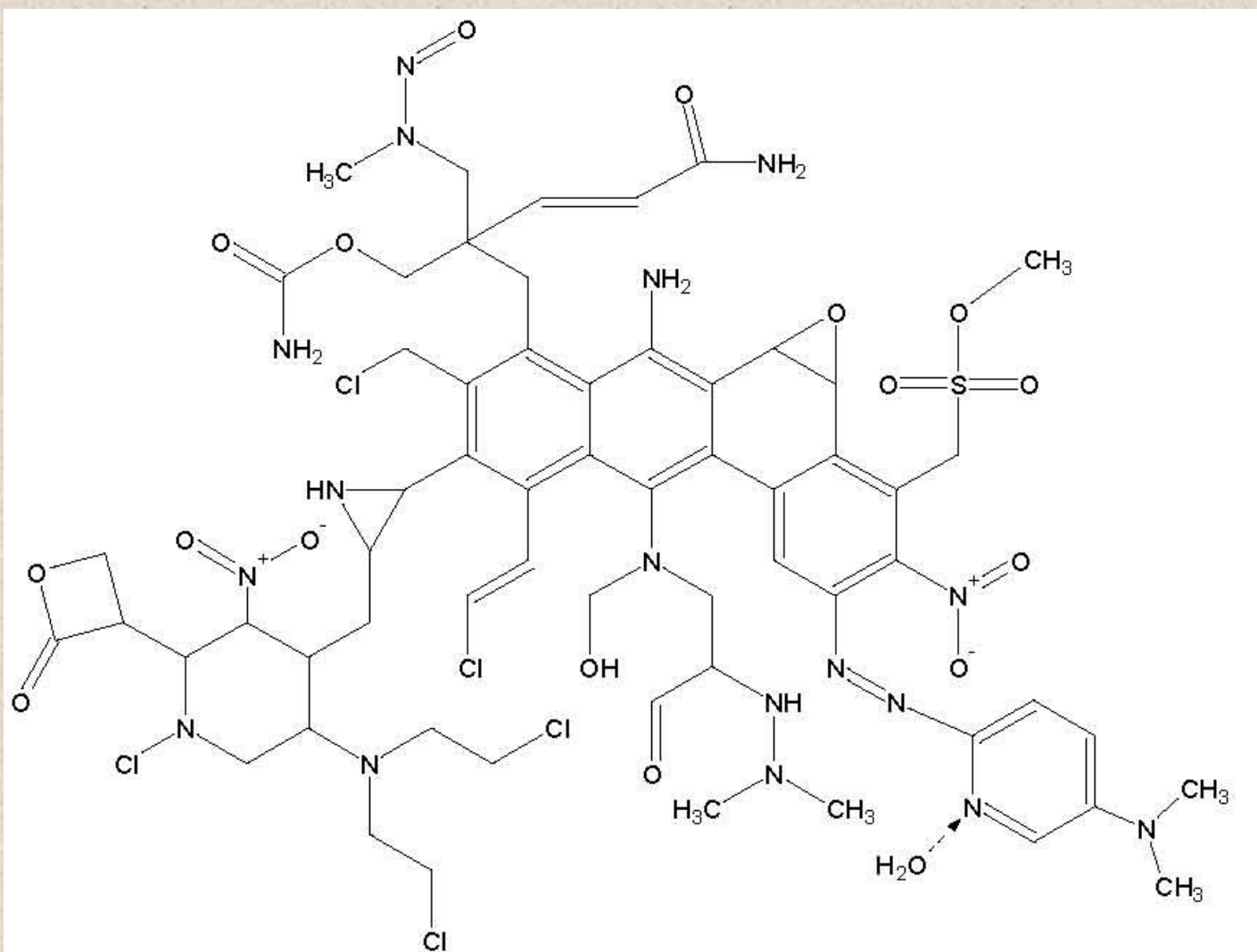
$\varepsilon_{\text{HВМО}}$ Енергія найнижчої вакантної МО (eV). розрахунок AM1

$\lg P$ Ліпофільність α - Поляризованість (\AA^3)

Ефективність (LOO) = 91 %

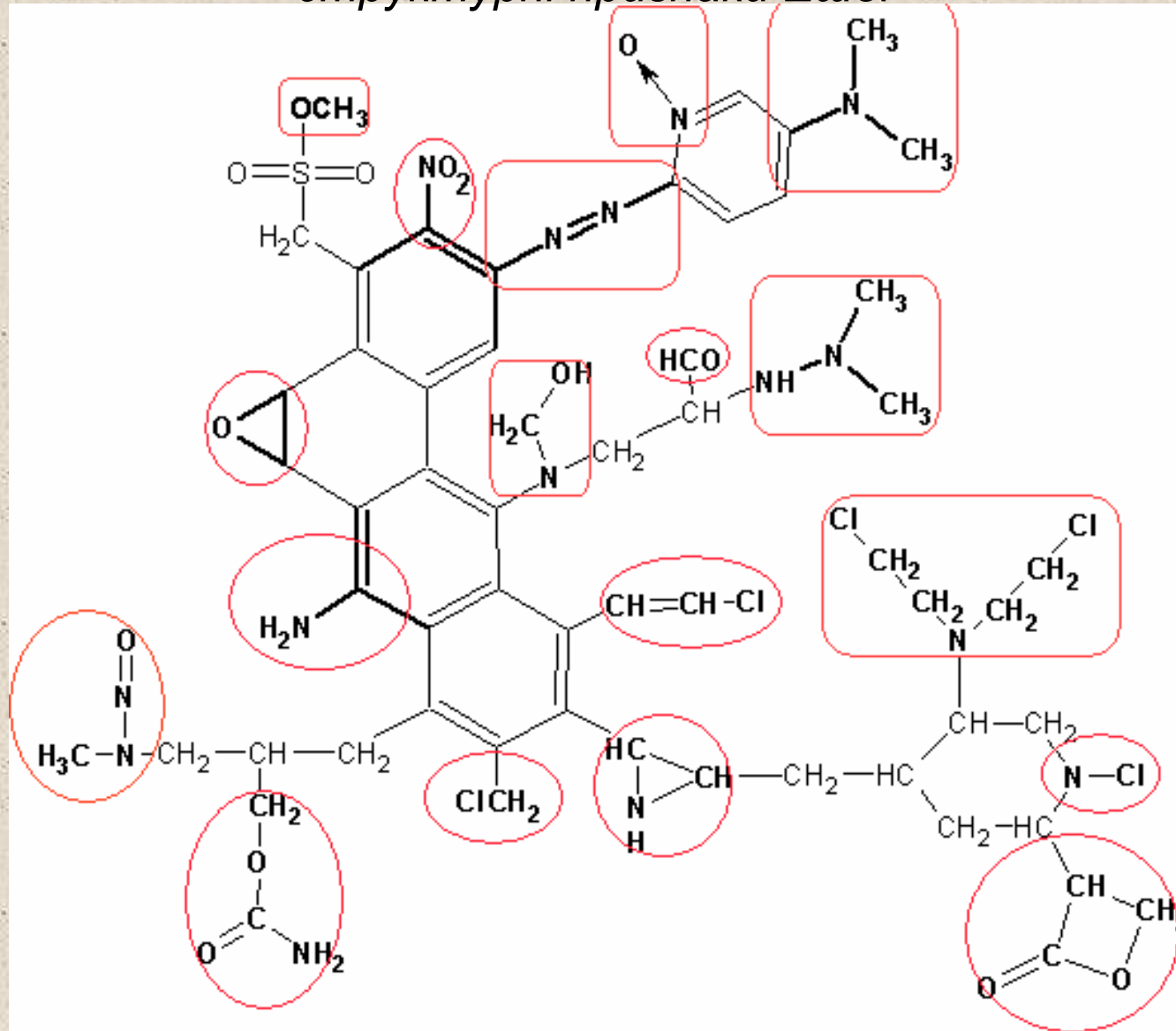
Модель поліканцерогена Ешбі (Ashby)

«структурні признаки Ешбі»



Модель поліканцерогена Ешбі (Ashby)

«структурні признаки Ешбі»



Розповсюдження РАН

Звідки беруться РАН в довколишньому середовищі ?

- Неповне сгоряння палива
- Дизельні мотори
- Індустрія
- Згоряння дерева

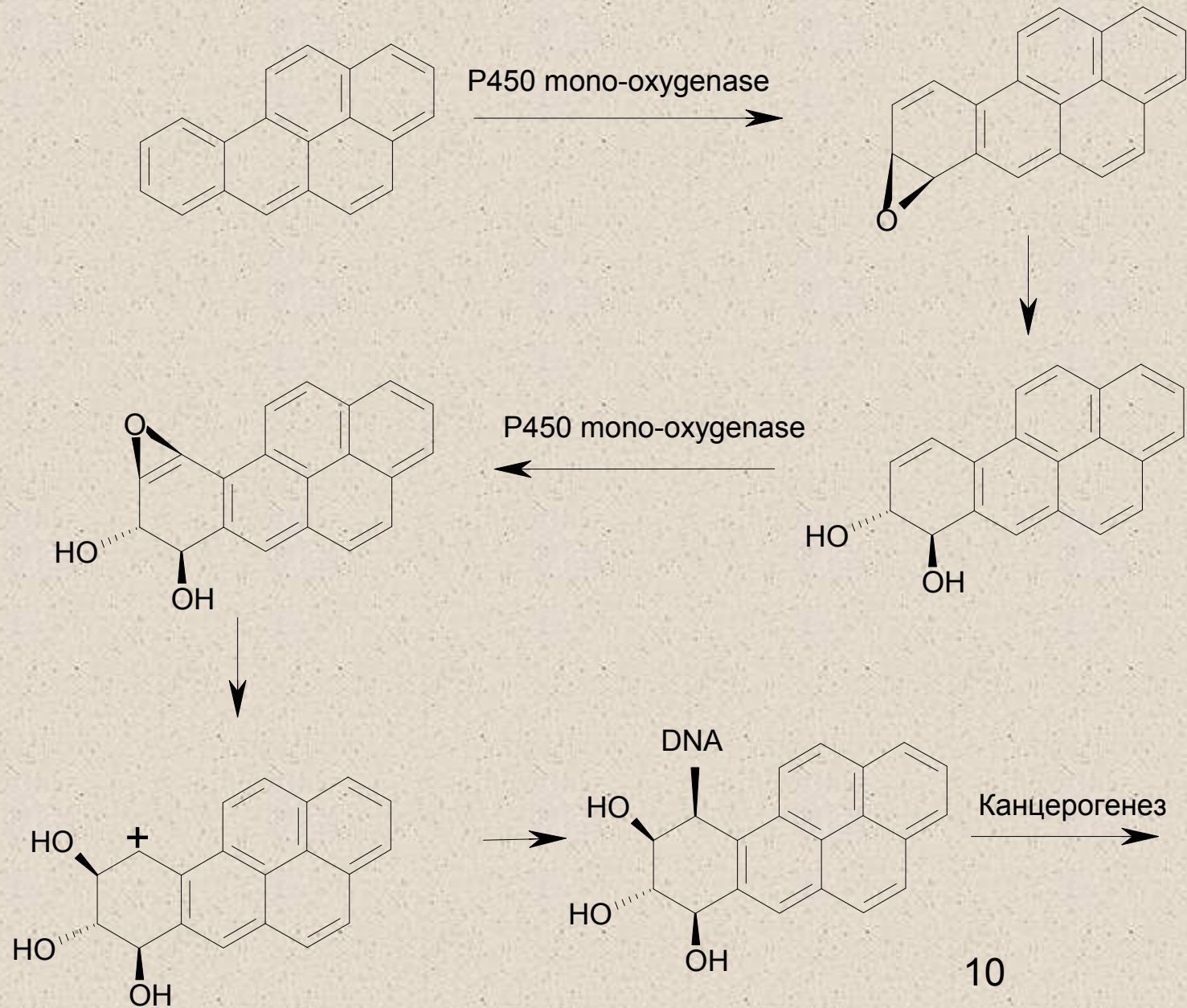
Як потрапляють РАН в організм людини ?

- Їжа
- Повітря

Згідно Стокгольмської конвенції РАН є постійними органічними забруднювачами:

- **Тератогенність**
- **Мутагенність**
- **Канцерогенність**

Метаболізм РАН в організмі

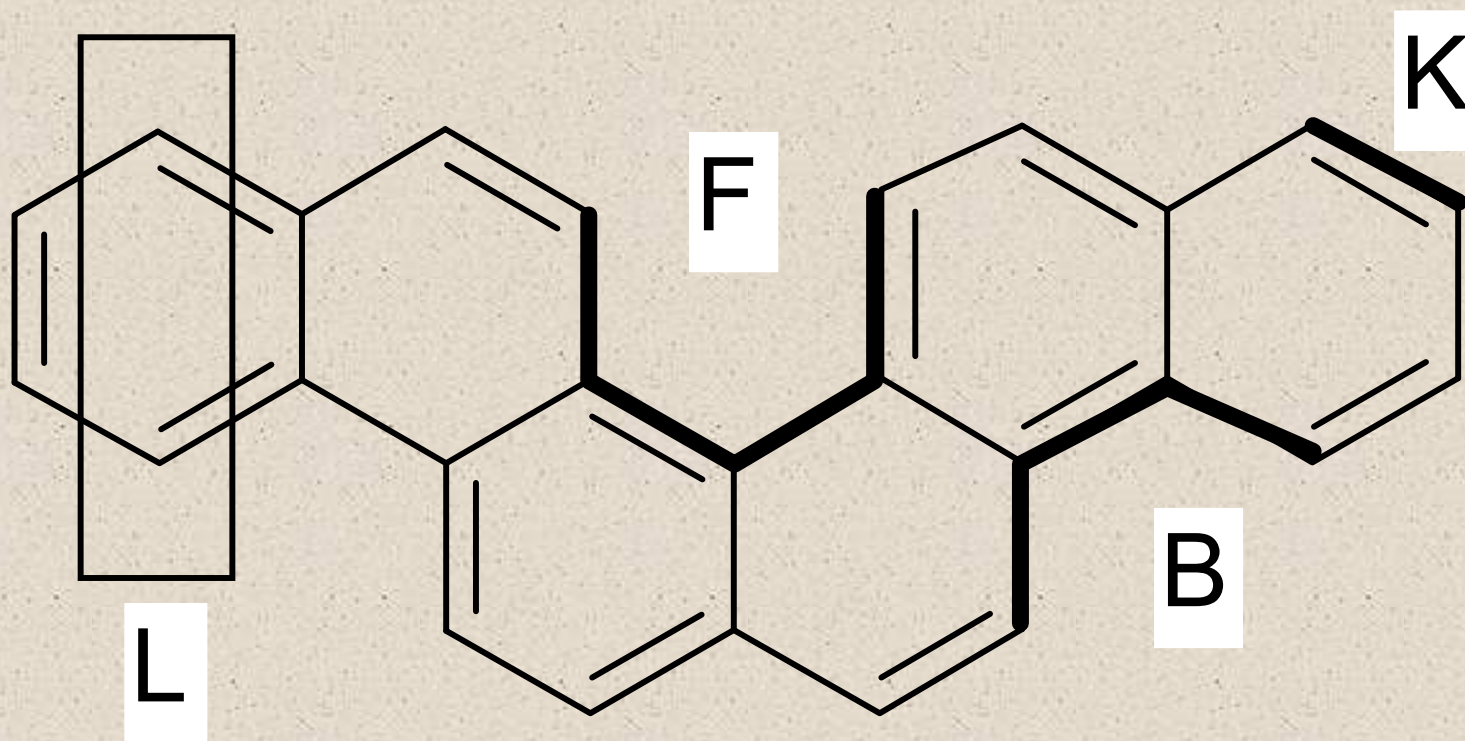


Цитохром Р450

For Educational Use Only

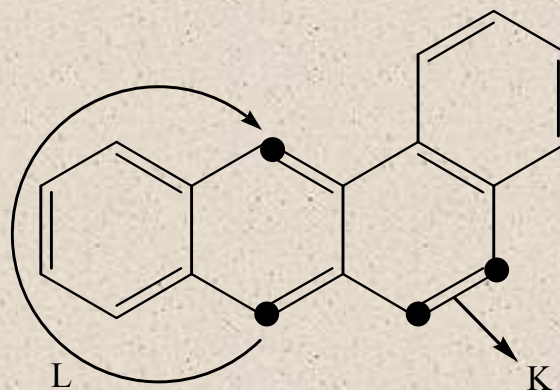


Активні області РАН



Дискримінантний аналіз

Канцерогенна активність конденсованих вуглеводнів

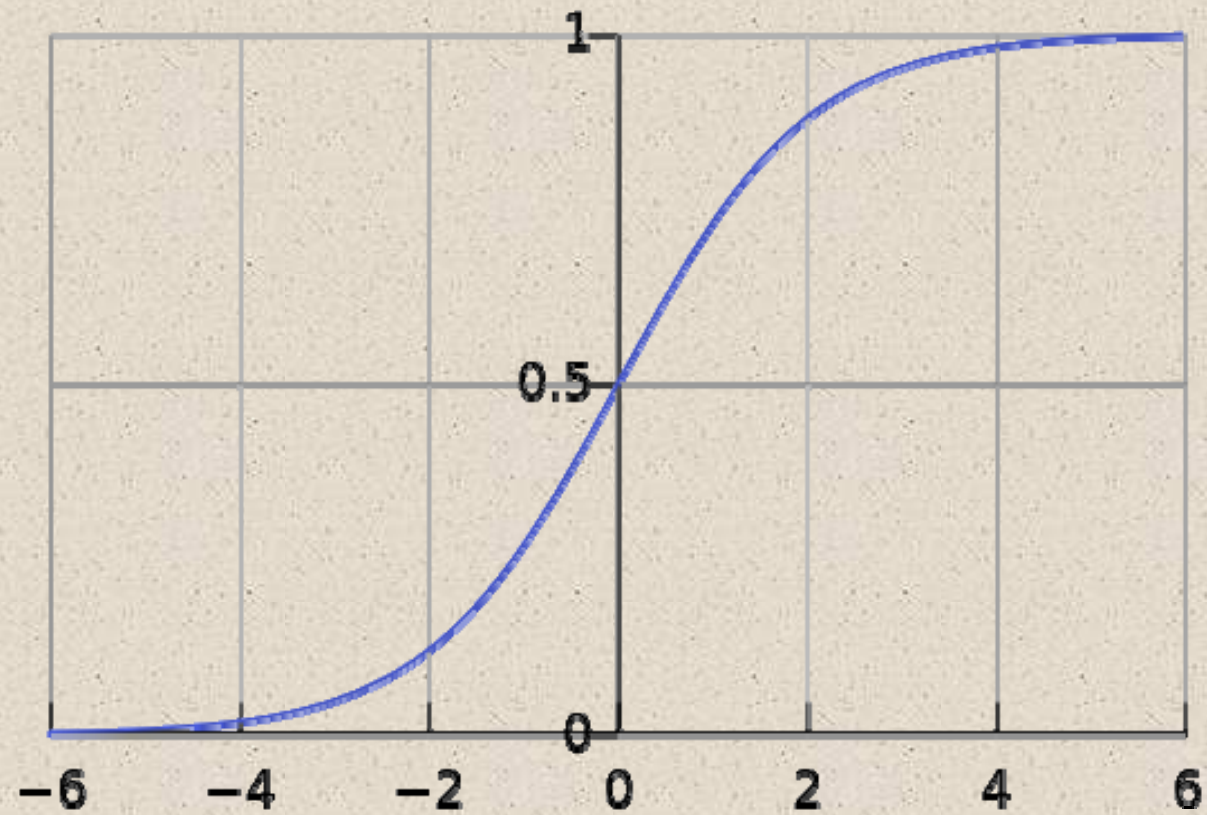


$$-24.1 + 17.7S_e(K) - 4.2S_e(L) > 0$$

S_e – “суперделокалізованість” (індекс Фукуї) K і L області.

Местечкин М.М., Сивякова, Канцерогенная активность полициклических Углеводородов (препринт Института Углекими, Донецк.)¹³

Логістична регресія



Логістична регресія

- **Логістична модель** найбільш популярна бінарна модель

- **Логістична модель**

Відгук системи:

$Y = 1$ (істина, да, активно)

$Y = 0$ (не істина, ні, неактивно)

В загальному випадку $Y = [0 - 1]$

- **Логістична модель** може бути узагальнена на випадок більше ніж два відгуки (*multinomial logistic regression model*)

Логістична регресія

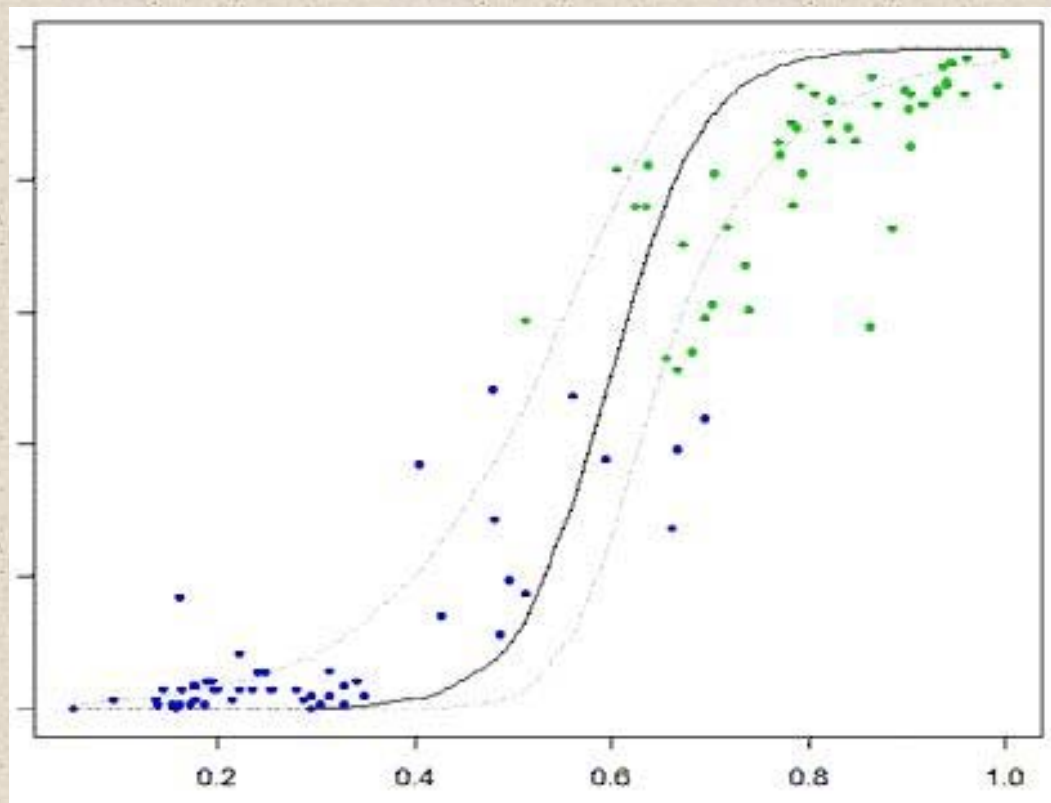
$$y_j = \frac{\exp(t_j)}{1 + \exp(t_j)} = \frac{1}{1 + \exp(-t_j)}$$

$$t_j = a_0 + \sum_i a_i d_{ji}$$

d_{ji} – дескриптори
(j -номер молекули,
 i - номер дескриптору)

$$\pi = Y_{(\text{calc})}$$

відгук системи
(клас)



максимізація $L = \prod_{j=1}^N \pi_j^{y_j} (1 - \pi_j)^{1-y_j}$

$$\ln(L) = \sum_{j=1}^n \{y_j \ln [\pi_j] + (1 - y_j) \ln [1 - \pi_j]\} \quad 16$$

Логістична регресія

$$L = \prod_{j=1}^N \pi_j^{y_j} (1 - \pi_j)^{1 - y_j} \quad \pi_j - \text{розраховані значення } y_j$$

π	y	$\pi^y (1 - \pi)^{(1-y)}$	L	$\ln(L)$
1	1	$1^1 (1-1)^{(1-1)}=1$	$1 \cdot 1=1$	0
0	0	$0^0(0-0)^{(0-0)}=1$		
$\frac{1}{2}$	1	$(1/2)^1 (1-1/2)^0=0.5$	$0.5 \cdot 0.5 = 0.25$	-1.386
$\frac{1}{2}$	0	$(1/2)^0 (1-1/2)^1=0.5$		
0	1	$0^1 (1-0)^{(1-1)}=0$	$0 \cdot 0 = 0$	$-\infty$
1	0	$1^0 (1-1)^{(1-0)}=0$		

Якість апроксимації

1. Процент вірно класифікованих систем, η
2. Логарифм функції правдоподібності, $\ln(L)$
3. Інформаційний критерій, $AIC = -2(\ln(L) - 2)$
4. Коефіцієнт детермінації, R^2
5. Урегульований (adjusted), R^2

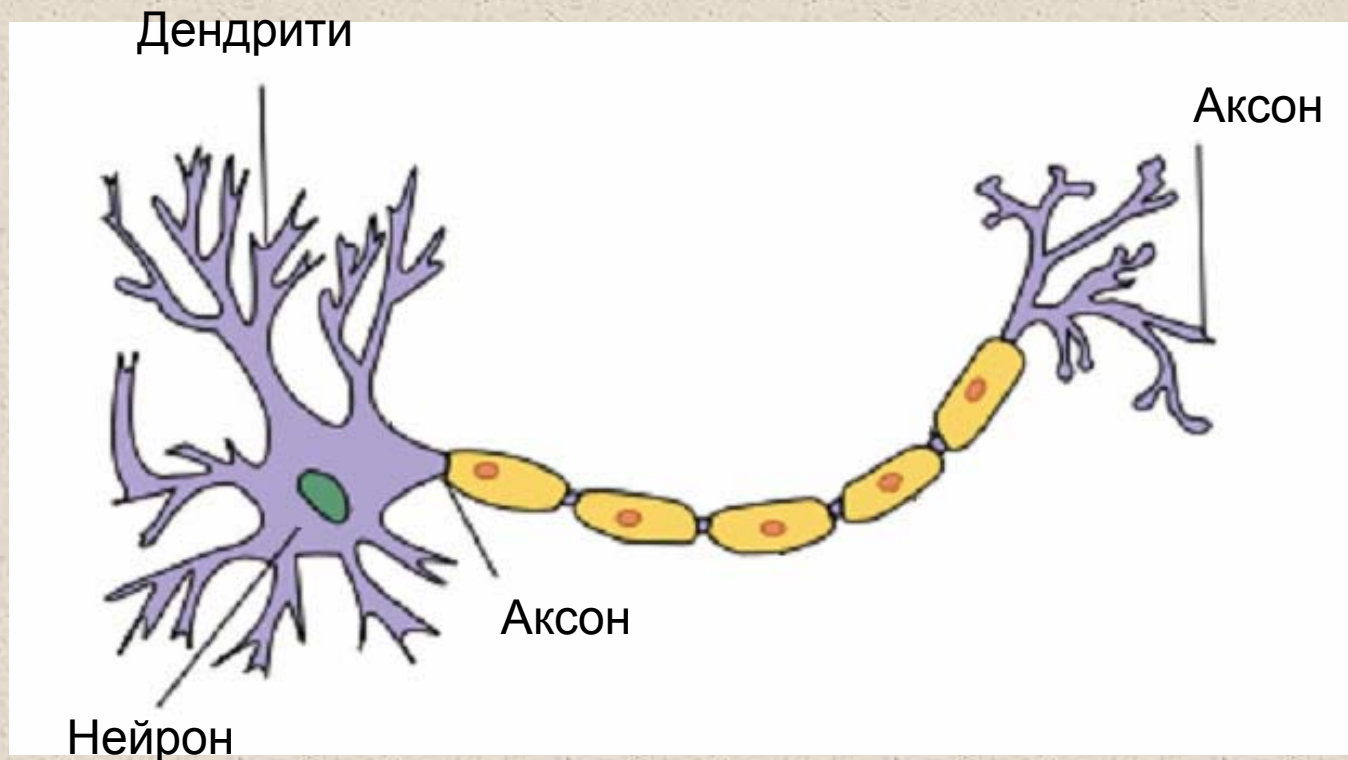
Процедура LOO

N

d_1	d_2	d_3	d_p	Class
					1
					0
					...

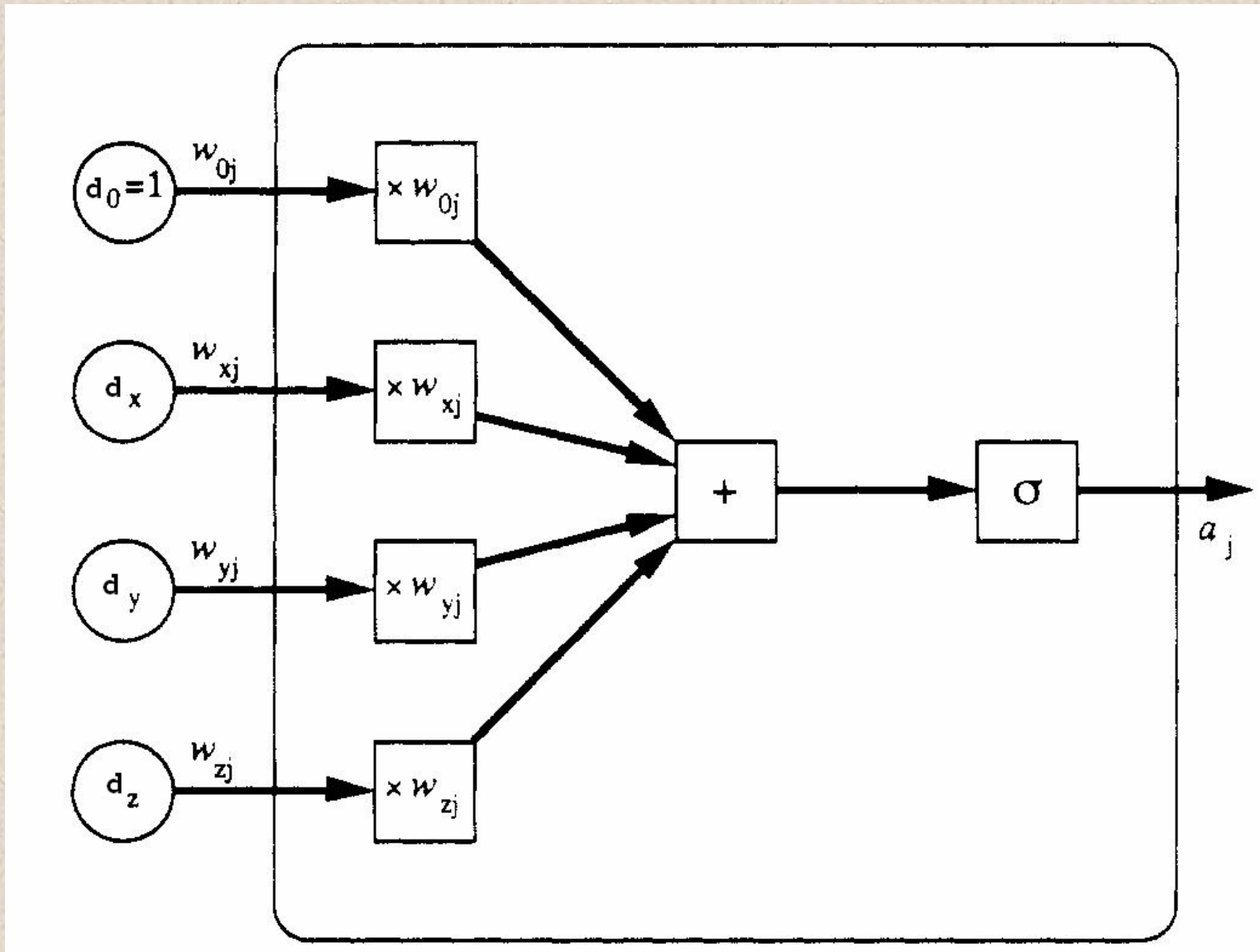
Метод Нейронных сетей

Структура нейрона

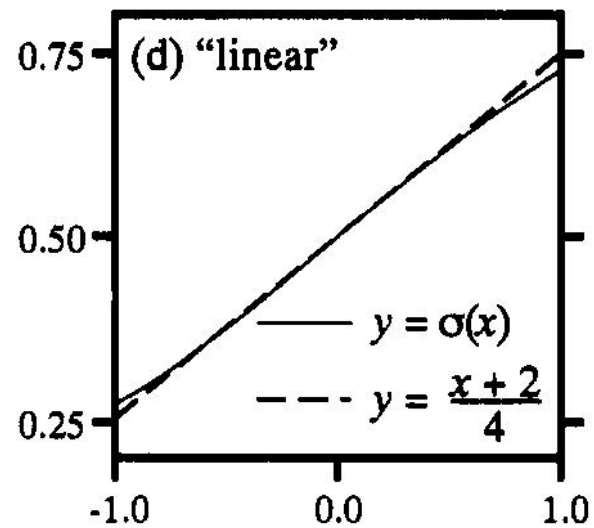
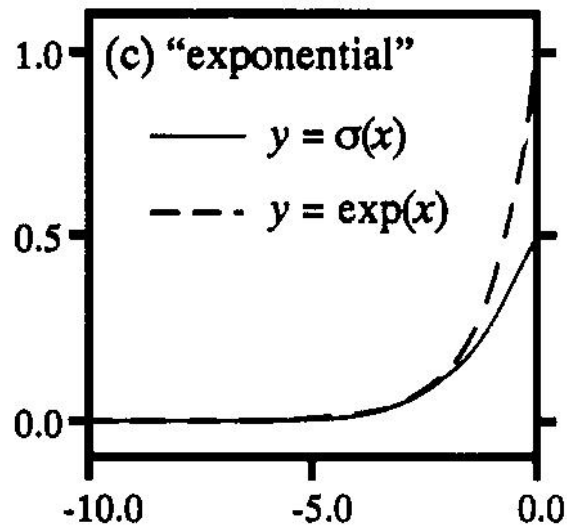
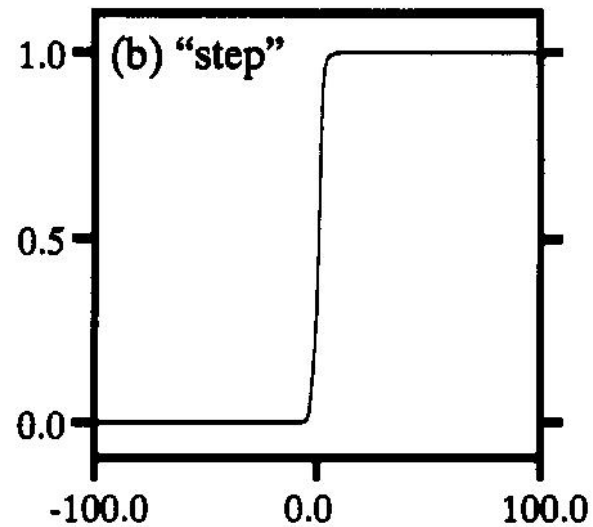
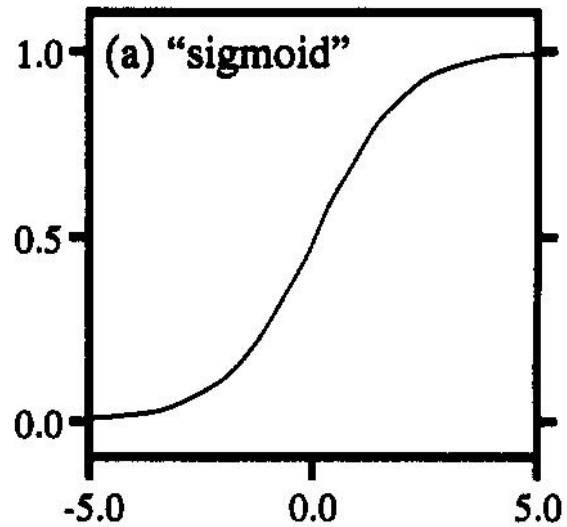


В человеческом организме $\sim 10^{15}$ нейронов

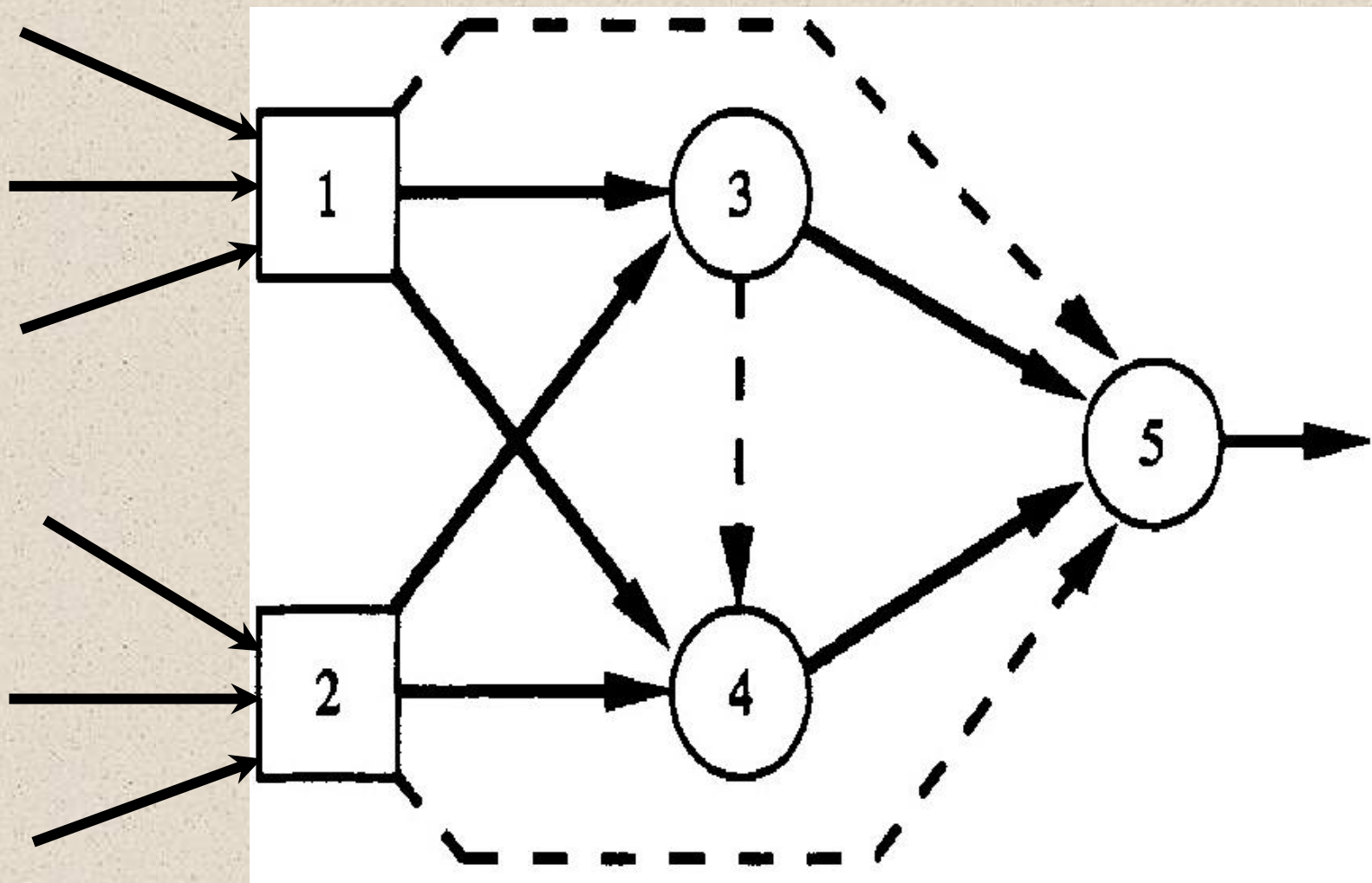
Модель Нейрона



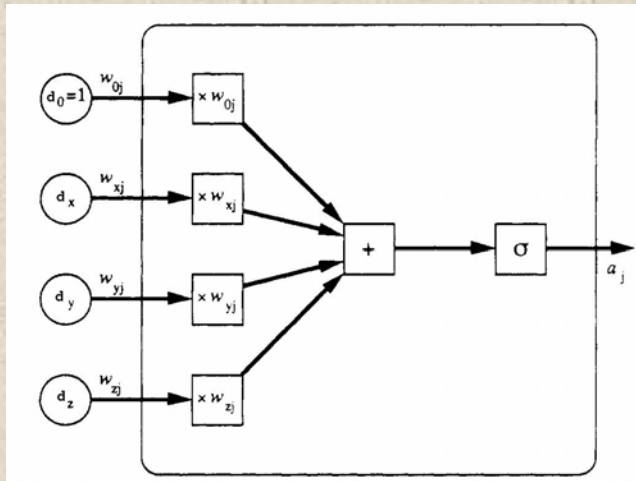
Вихідна функція σ



Модель Нейроної мережі



Навчання Нейроної мережі



Начальне наближенн W

Вычисляем a_j

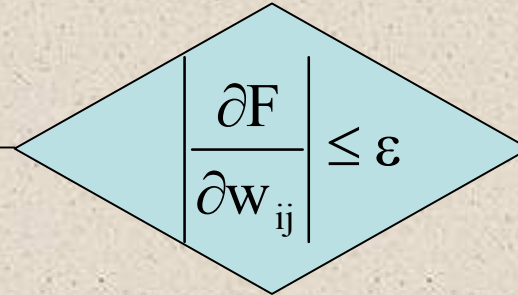
Вычисляем
Функцию ошибки F

d_0	d_x	d_y	d_z	a
				a_1
				a_2
				a_3
				a_4

$$F = \sum_i (a_i - a_i^{(calc)})^2$$

нет

$$W_{ij}^{(n+1)} = W_{ij}^{(n)} - \zeta \frac{\partial F}{\partial w_{ij}}$$



да

STOP