
GMDH Applications in QSAR and QSAR ensembling

Giuseppina Gini

Department of Electronics and Information,
Politecnico di Milano, Italy



What is QSAR?

- ◆ A QSAR is a **mathematical relationship (model)** between a biological activity of a molecular system and its geometric and chemical characteristics.
- ◆ The problem of QSAR is to find coefficients C_0, C_1, \dots, C_n such that the obtained linear (or non linear) equation:

$$\text{Biological activity} = C_0 + (C_1 * P_1) + \dots + (C_n * P_n)$$

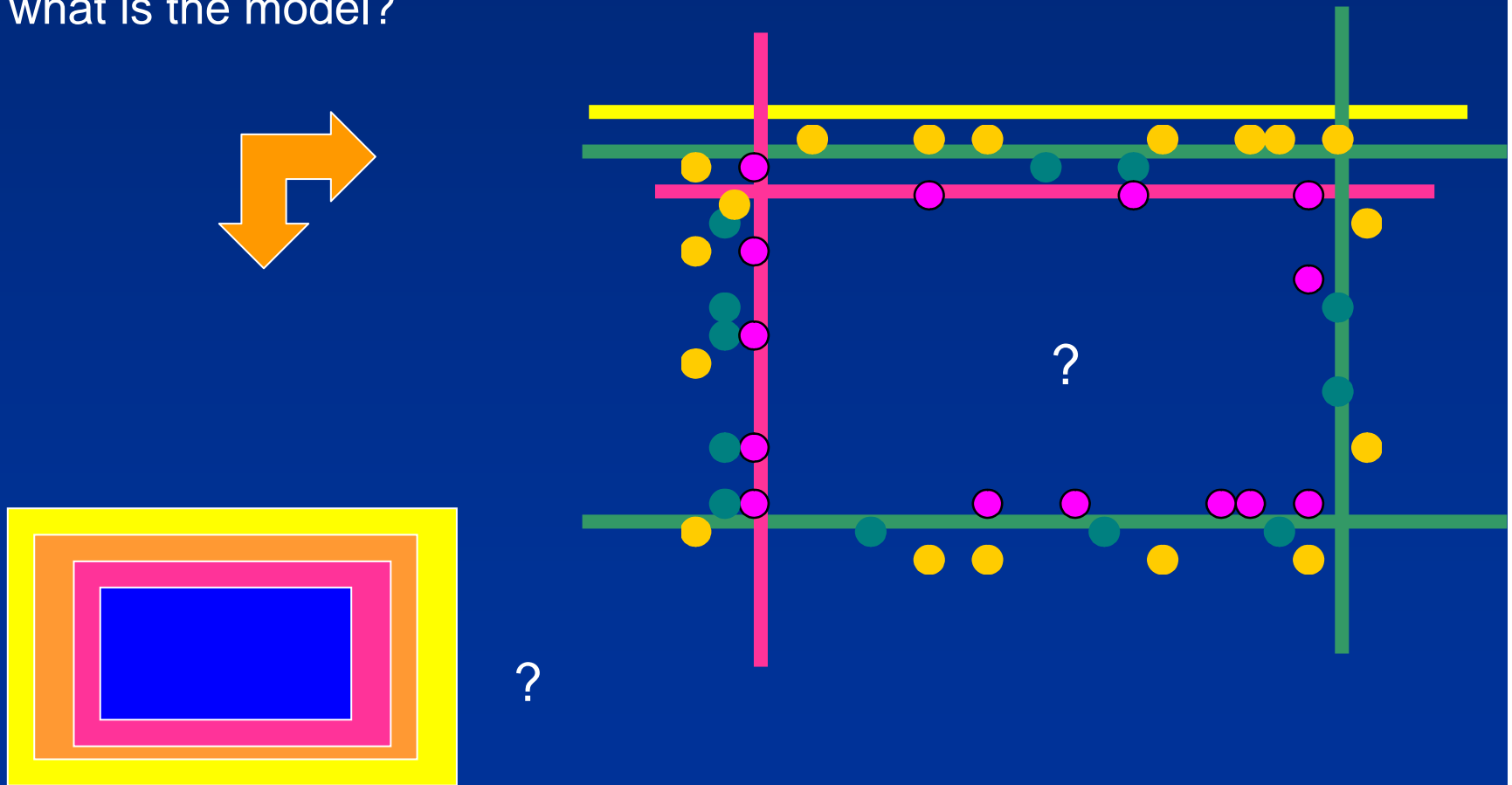
minimizes the prediction error.

Also QSAR attempts to find consistent relationship between biological activity and molecular properties, so that these “rules” can be used to evaluate the activity of new compounds.

What is a model?

Inducing a model

- Given the illustrated points, what is the model?



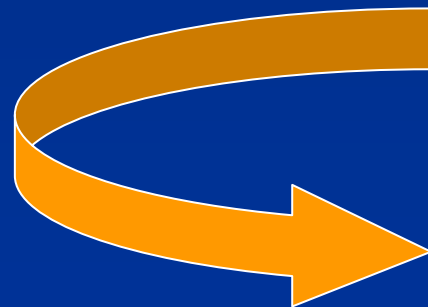
Trends

- from Statistics to Intelligent Data Analysis to Machine Learning
 - move from basic statistics to more complex models
 - Lack of adaptation
 - Noisy data
 - Computational problem Trend from exact to approximate solutions

Exact solution

optimization

NP hard



Approximate solution

"Soft computing"

Modelling: scientific communities

- statistics
- machine learning
- pattern recognition
- neural networks
- computational learning theory
- knowledge discovery and data mining
- ...

The equivalence of algorithms

- Artificial neural networks
- Evolutionary Algorithms
- Bayesian networks
- Decision trees
- ...



More interpretability?

EQUIVALENCE - if algorithm A outperforms algorithm B on some learning task, then there must exist exactly as many other tasks where B outperforms A.

Techniques in QSAR (ISI 1990-2005)

Query in text	Number of items
QSAR or predictive toxicology	4835
QSAR (in text)	4738
QSAR (in title)	1860
<i>Technique used</i>	
QSAR and regression	950
QSAR and NN	237
QSAR and classification	223
QSAR and PCA	104
QSAR and GA	95
QSAR and ML	20
QSAR and ensemble	19
QSAR and SVM	17
QSAR and Radial	14
QSAR and fuzzy	10
QSAR and Bayesian	4

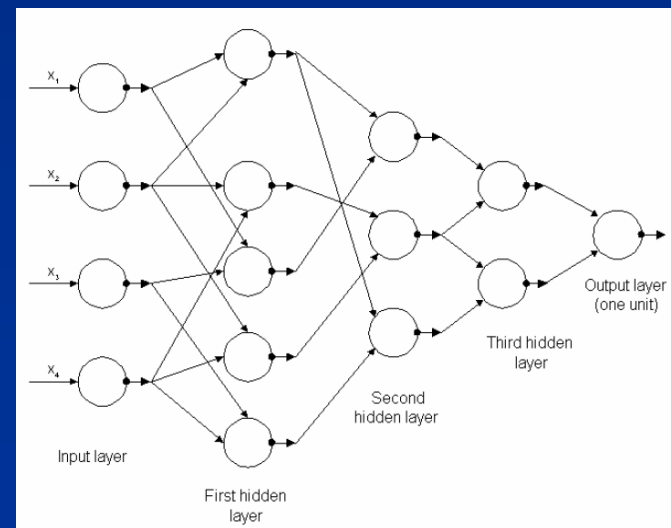
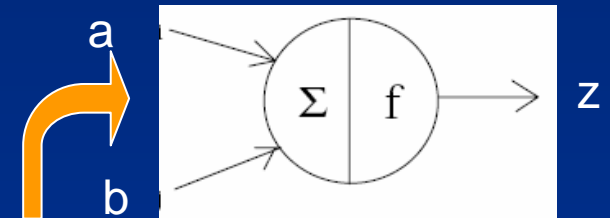
Proposal of the modelling method: GMDH

- *GMDH (Group Method of Data Handling, A.G. Ivakhnenko, 1967)*
- Proposed both for the QSAR activity and for the integration of different models (ensembling)
- Is auto-organizing
- Uses and external criterion
- Avoids overfitting of data
- Does not require expertise in ANN construction

Description (1)

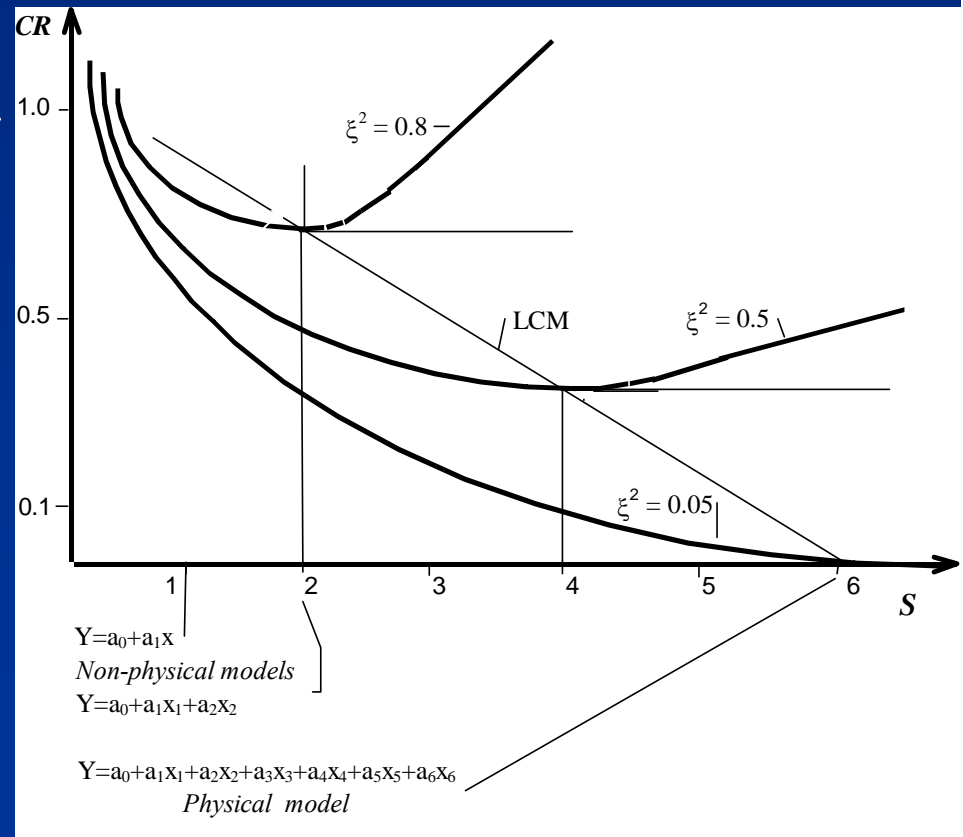
1. automatically creates the network structure
2. neurons are 2-levels structures; nets are multi-layered
3. learning is contextual to net construction; w are found by solving Linear Regression equations with $z = y$
4. the structure is designed by gradually increasing complexity
5. it avoids overfitting through a regularity criterion applied to a validation subset of the training set
6. pruning criteria use the results on the validation subset

$$z = w_0 + w_1a + w_2b + w_3a^2 + w_4b^2 + w_5ab$$

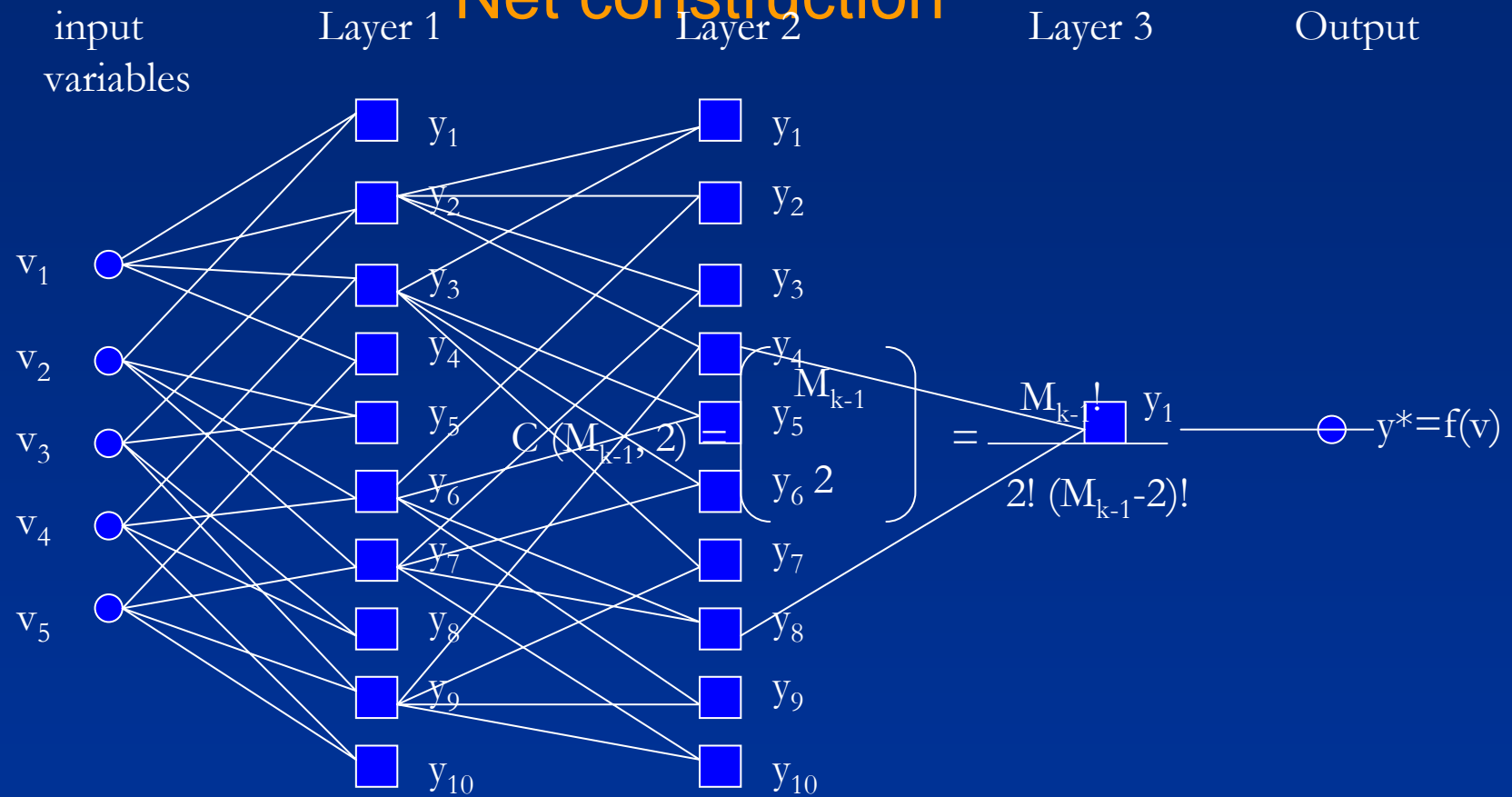


Description (2)

- Inductive algorithms select the optimal non physical model
- We see the minimum values of the external criterion (LCM) versus the complexity of the structure S of the model for different values of noise variance ξ^2 .



Net construction



...our algorithm: poliGMDH

- Multilayer Algorithm
- Ivakhnenko polynomial

$$z^k = a(z_i^{k-1})^2 + b(z_i^{k-1})(z_j^{k-1}) + c(z_j^{k-1})^2 + d(z_i^{k-1}) + e(z_j^{k-1}) + f$$

- In each layer $C(n_{(k-1)}, 2)$ neurons
- Selection of good neurons
 1. m.s.e. value on a selection set
 2. Neurons are ordered
 3. The last are eliminated

$$C(n_{(k-1)}, 2) = \binom{n_{k-1}}{2}$$

Polynomials

- find the polynomial of optimal complexity
- high order polynomials cause overfitting.
- second order polynomials are preferred

$$p_1(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2$$

$$p_2(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2$$

$$p_3(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2$$

$$p_4(\mathbf{x}) = w_0 + w_1x_1 + w_2x_1x_2 + w_3x_1^2$$

$$p_5(\mathbf{x}) = w_0 + w_1x_1 + w_2x_1x_2$$

$$p_6(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1^2$$

$$p_7(\mathbf{x}) = w_0 + w_1x_1 + w_2x_1^2 + w_3x_2^2$$

$$p_8(\mathbf{x}) = w_0 + w_1x_1^2 + w_2x_2^2$$

$$p_9(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2 + w_4x_1^2 + w_5x_2^2$$

$$p_{10}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2 + w_4x_1^2$$

$$p_{11}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_1x_2 + w_3x_1^2 + w_4x_2^2$$

$$p_{12}(\mathbf{x}) = w_0 + w_1x_1x_2 + w_2x_1^2 + w_3x_2^2$$

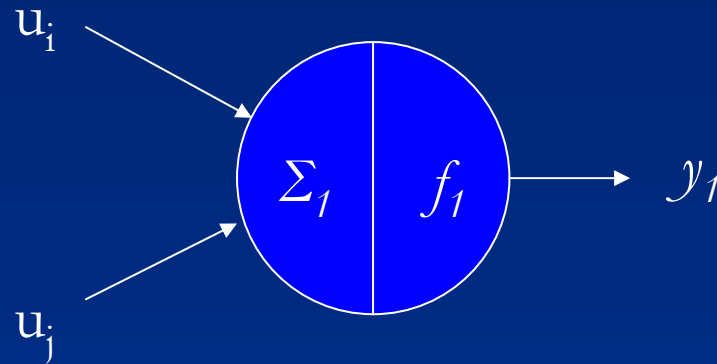
$$p_{13}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_1x_2 + w_3x_2^2$$

$$p_{14}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2^2$$

$$p_{15}(\mathbf{x}) = w_0 + w_1x_1x_2$$

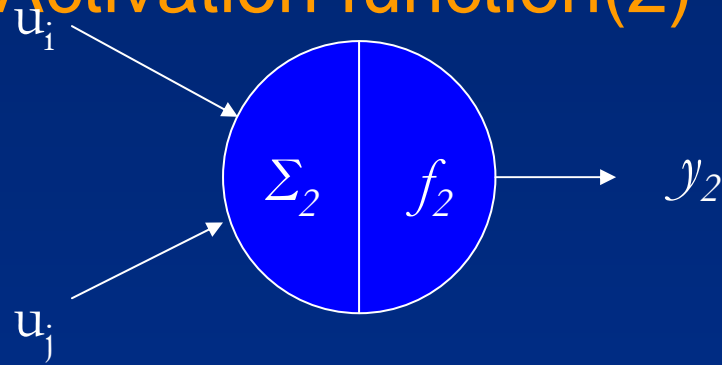
$$p_{16}(\mathbf{x}) = w_0 + w_1x_1x_2 + w_2x_1^2$$

Activation function (1)



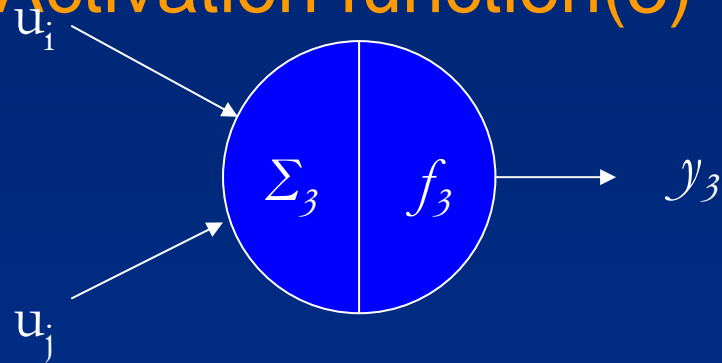
- Ivakhnenko:
- Σ_1 : (*non-linear*)
- $z_1 = A * x_1 * x_2 + B * x_1^2 + C * x_2^2 + D * x_1 + E * x_2 + F$
- f_1 : (*linear*)
- $y_1 = z_1$

Activation function(2)



- bi-quadratic:
- Σ_2 : (*non-linear*)
- $z_2 = A * x_1 * x_2 + B * x_1 + C * x_2 + D$
- f_2 : (*linear*)
- $y_2 = z_2$

Activation function(3)



- Exponential
- Σ_3 : (non-linear)
- $z_3 = y_2 = z_2 = A * x_1 * x_2 + B * x_1 + C * x_2 + D$
- f_3 : (non-linear)
- $y_3 = \exp(-z_3^2) = \exp(-y_2^2)$

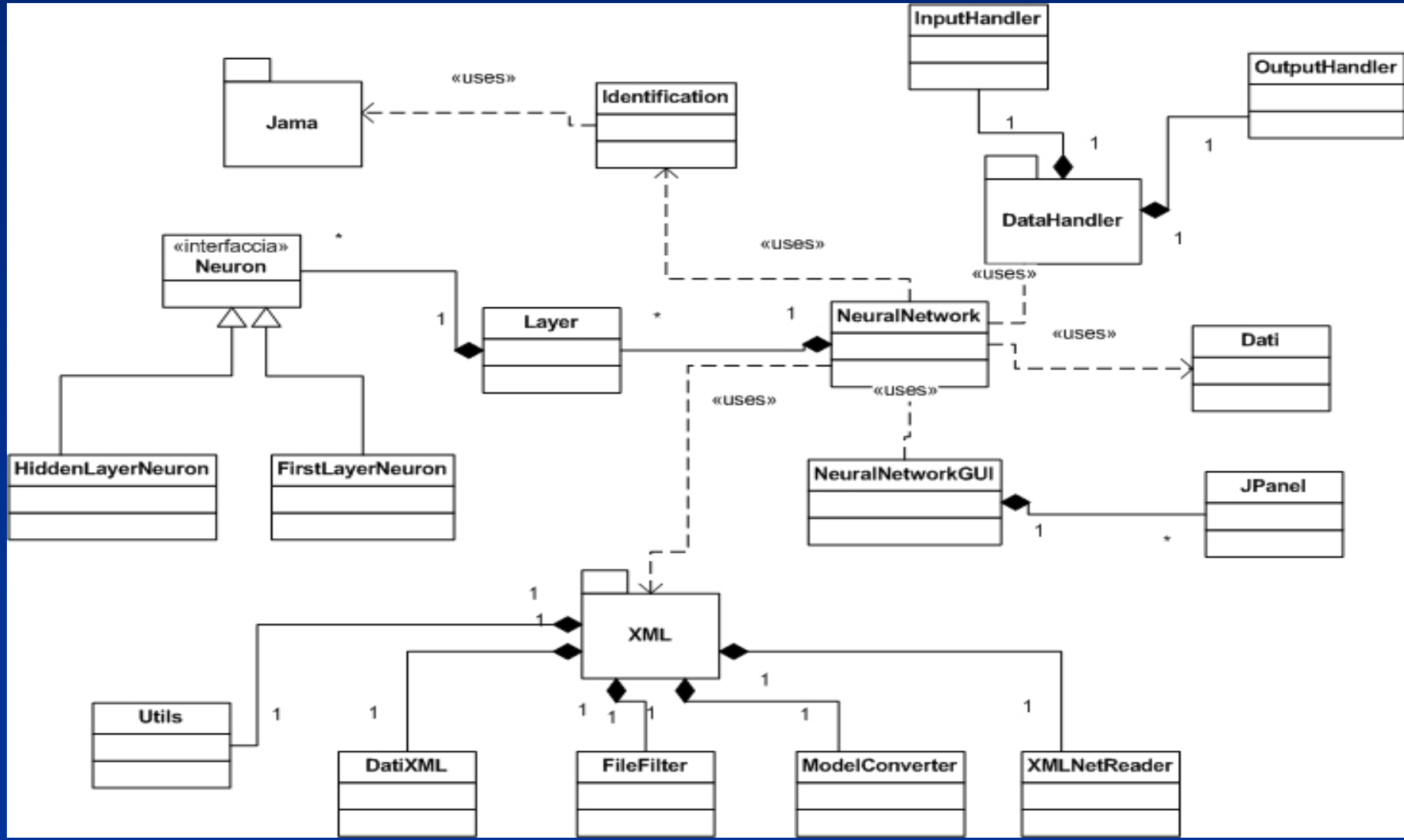
Transfer functions

- 1) Ivakhnenko:
 - Σ_1 : (*non-linear*)
 - $z_1 = A*x_1*x_2 + B*x_1^2 + C*x_2^2 + D*x_1 + E*x_2 + F$
 - f_1 : (*linear*) $y_1 = z_1$

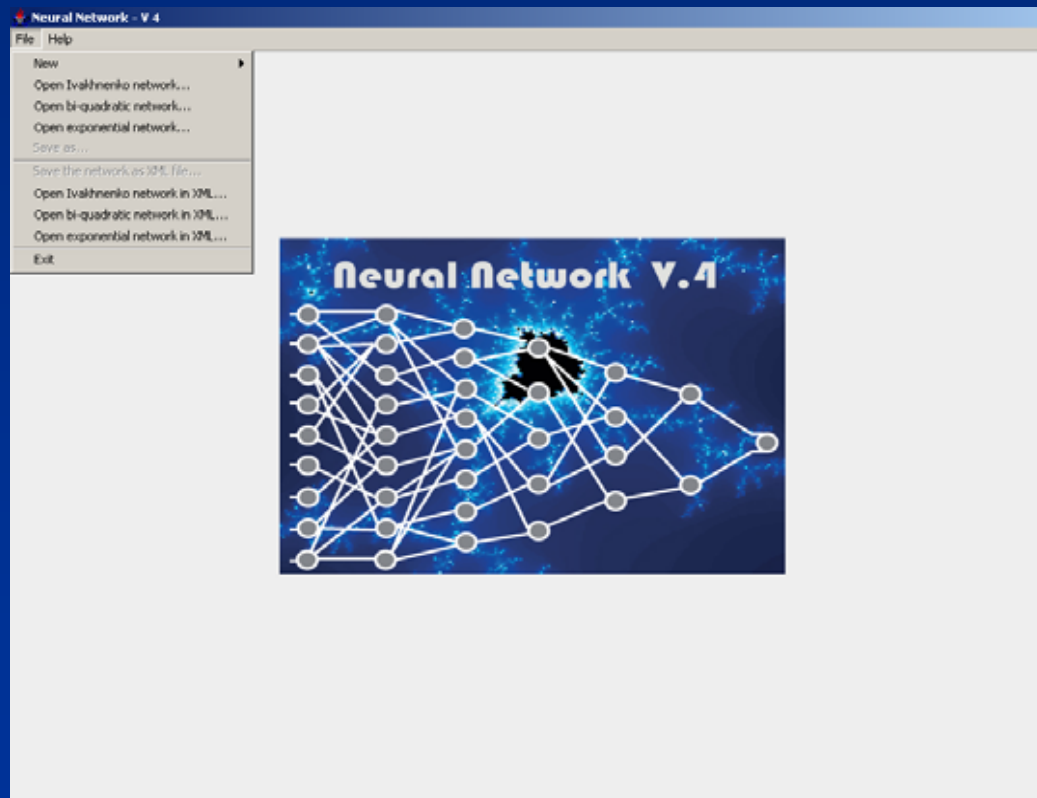
- 2) bi-quadratic:
 - Σ_2 : (*non-linear*)
 - $z_2 = A*x_1*x_2 + B*x_1 + C*x_2 + D$
 - f_2 : (*linear*) $y_2 = z_2$

- 3) exponential:
 - Σ_3 : (*non-linear*)
 - $z_3 = y_2 = z_2 = A*x_1*x_2 + B*x_1 + C*x_2 + D$
 - f_3 : (*non-linear*) $y_3 = \exp(-z_3^2) = \exp(-y_2^2)$

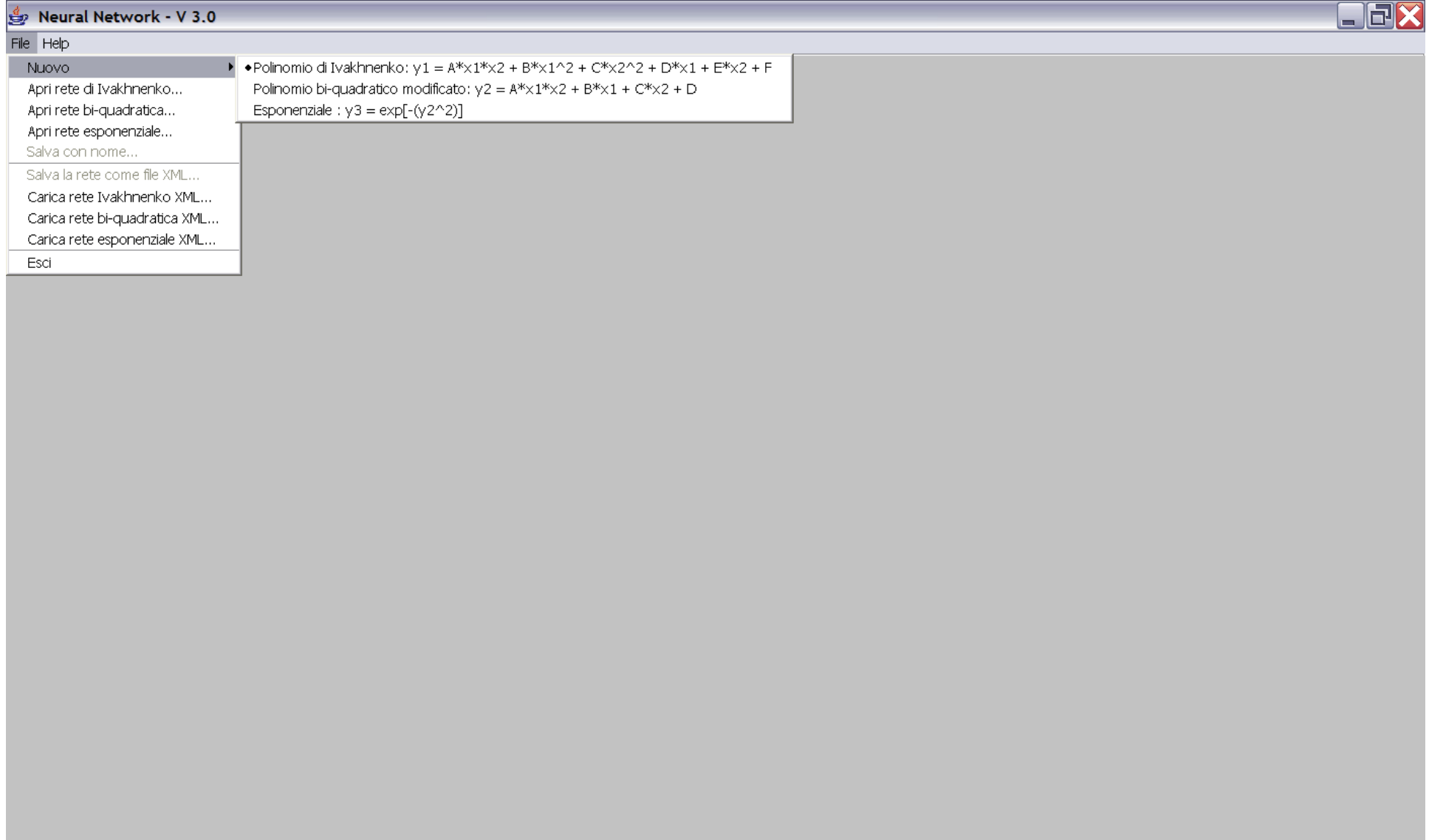
software – Class Diagram



The Neural Network construction in poliGMDH



File menu



Building the net

Scelta parametri di costruzione

Selezionare il file di dati:

Numero massimo per strato nella rete:

Percentuale di dati da usare come train set:

Scelta parametri di costruzione

Dati usati per la creazione della rete

Scegliere l'ordinamento per i dati e procedere

	INO	IN1	IN2	IN3	IN4	IN5	IN6	IN7	IN8	IN9	OUT
158.27	29.0	28.0	44.877	17.07	1509.0	25.632	76.0	700.0	2520.0	-0.55344...	
238.46	47.0	46.0	77.48	17.07	4690.0	26.603	234.0	1004.0	3512.0	0.294119...	
197.03	17.0	17.0	39.852	59.92	1348.642	54.683	284.375	2590.0	14512.0	-5.37766...	
110.97	9.0	8.0	25.277	0.0	54.309	4.243	21.0	164.0	488.0	-3.33807...	
166.22	24.0	25.0	53.816	17.07	3647.0	16.793	3596.623	2438.0	12770.0	-1.89472...	
221.04	19.0	19.0	46.865	26.3	3136.309	43.186	572.563	1830.0	8816.0	-0.69786...	
335.25	41.0	41.0	60.909	44.76	8036.751	263.482	1107.813	2618.0	12102.0	-5.86616...	
321.22	38.0	38.0	75.551	44.76	9196.196	131.662	1037.406	2376.0	10950.0	-5.03496...	
277.16	31.0	31.0	65.507	35.53	5220.309	79.013	632.166	2166.0	10194.0	-6.54068...	
347.31	46.0	46.0	88.276	35.53	8673.31	651.487	1179.219	2988.0	14214.0	-3.78335...	
218.69	26.0	27.0	62.294	0.0	3695.776	16.553	2748.576	2396.0	11802.0	-5.71815...	
141.61	17.0	17.0	40.604	0.0	607.333	11.269	272.531	1362.0	6678.0	-2.90999...	
101.17	18.0	18.0	25.133	9.23	274.0	5.568	46.0	1086.0	5134.0	-0.06082...	
282.26	33.0	33.0	71.12	55.26	2945.518	92.59	406.0	2086.0	10012.0	-11.5573...	
156.62	19.0	19.0	42.639	0.0	860.111	20.125	337.813	1700.0	8760.0	-6.07547...	
163.19	20.0	19.0	41.161	90.51	1026.259	38.626	90.0	1346.0	6498.0	-0.51004...	
269.8	38.0	38.0	74.075	26.3	4453.111	146.915	1178.547	2776.0	14536.0	-6.56526...	
259.6	36.0	36.0	73.928	26.3	4373.111	150.06	1259.969	2640.0	15046.0	-5.59768...	
190.3	26.0	25.0	47.399	75.99	1697.0	71.562	148.0	1398.0	6212.0	-5.82842...	
222.3	28.0	27.0	48.247	93.21000...	3901.0	424.108	242.0	2326.0	11956.0	-1.66635...	
227.38	32.0	32.0	69.568	63.97	2604.0	61.782	698.875	2042.0	9930.0	-4.26347...	
253.45	45.0	46.0	87.493	27.96	10877.0	265.716	14655.47...	3506.0	17246.0	-5.98281...	
275.5300...	21.0	22.0	68.241	38.67	4467.223	31.828	3227.623	2670.0	13334.0	-7.58480...	
548.7	79.0	80.0	135.737	137.1	57445.0	427563.9	9630.609	6224.0	31440.0	-17.8204...	
215.72	28.0	28.0	62.214	38.67	2433.222	51.74	605.125	1916.0	9268.0	-3.86990...	

casuale var OUT cresc var OUT decres Annulla Procedi

Using the net

Neural Network - V 3.0 - 1.gnn

File Help

Parametri rete | Andamento errore | Topologia rete | Utilizzo rete | Bontà rete | Funzione di regressione

Num Layer	Num Neur	Num Neurone	Num 1 padre	Num 2 padre	Param 0	Param 1	Param 2	Param 3	Param 4	Param 5
0	10	0	-	-	0.0	0.0	0.0	1.0	0.0	0.0
1	7	1	-	-	0.0	0.0	0.0	1.0	0.0	0.0
2	8	2	-	-	0.0	0.0	0.0	1.0	0.0	0.0
3	10	3	-	-	0.0	0.0	0.0	1.0	0.0	0.0
4	11	4	-	-	0.0	0.0	0.0	1.0	0.0	0.0
5	11	5	-	-	0.0	0.0	0.0	1.0	0.0	0.0
6	14	6	-	-	0.0	0.0	0.0	1.0	0.0	0.0
7	15	7	-	-	0.0	0.0	0.0	1.0	0.0	0.0
8	16	8	-	-	0.0	0.0	0.0	1.0	0.0	0.0
9	15	9	-	-	0.0	0.0	0.0	1.0	0.0	0.0
10	13									
11	13									
12	14									
13	16									
14	14									
15	12									
16	9									
17	8									
18	9									
19	12									
20	11									
21	10									
22	11									
23	14									
24	14									
25	14									
26	12									
27	12									
28	9									
29	8									
30	6									
31	3									
32	2									
33	1									

Neural Network - V 3.0 - 1.gnn

File Help

Parametri rete | Andamento errore | Topologia rete | Utilizzo rete | Bontà rete | Funzione di regressione

Usare il punto come separatore decimale

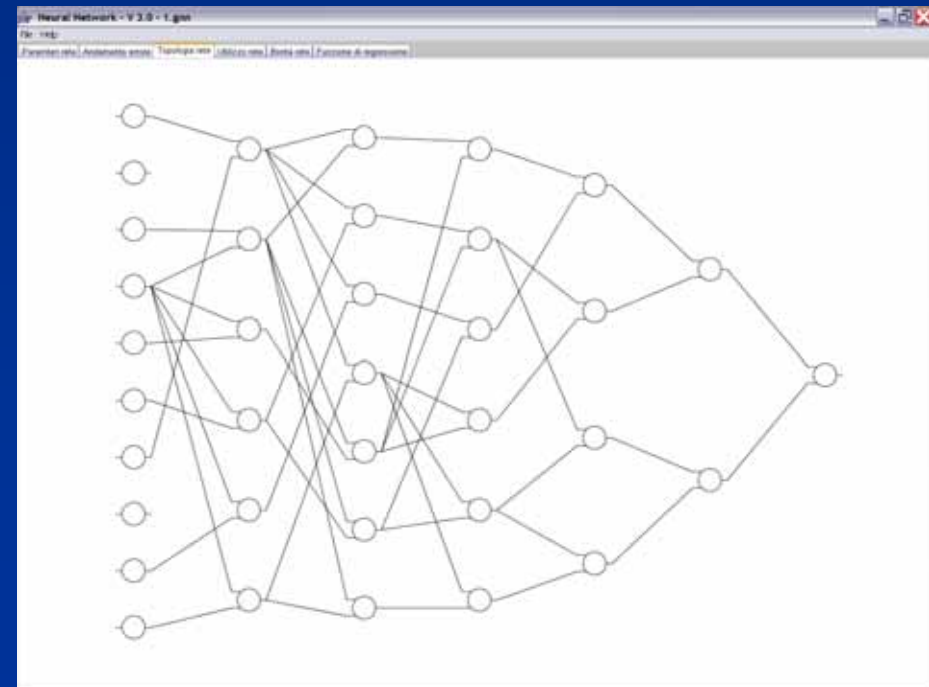
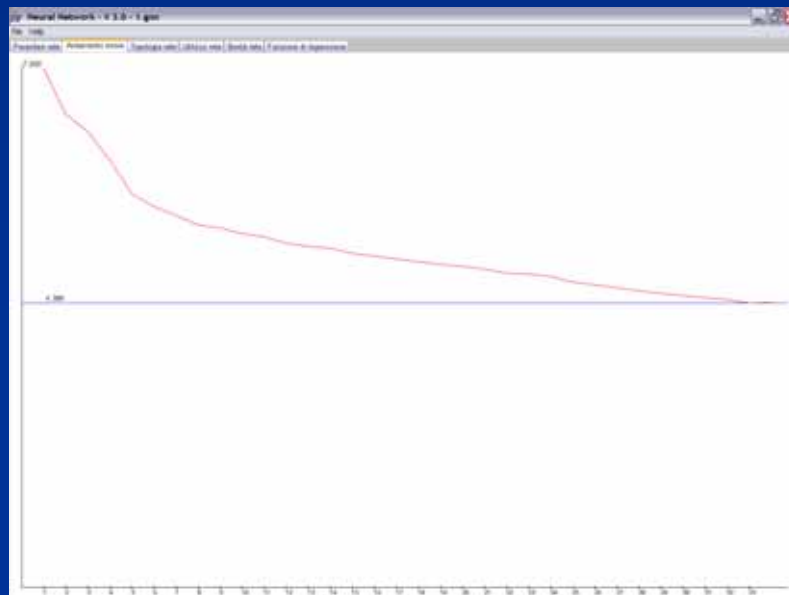
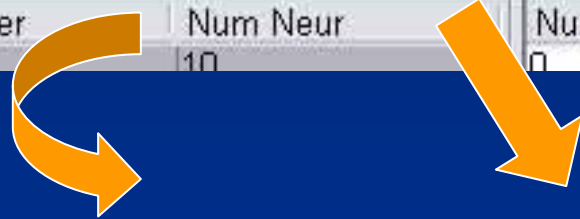
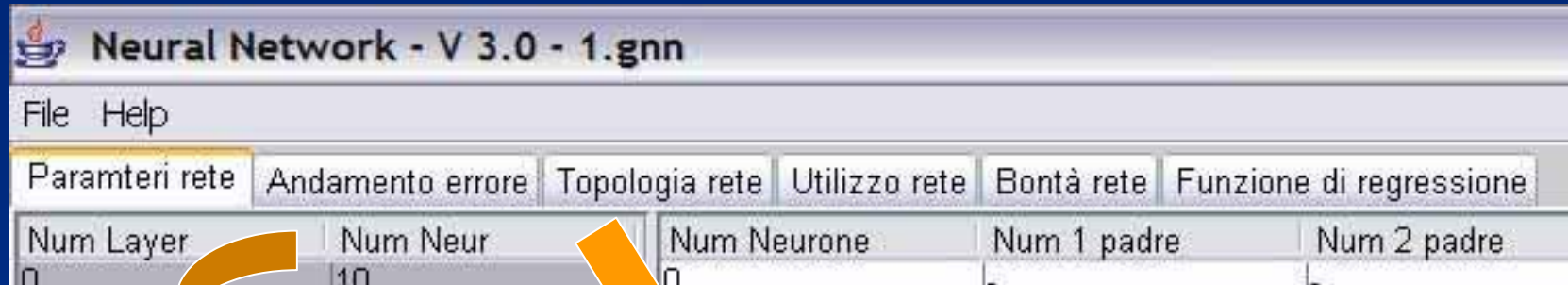
IN 1	0.1
IN 2	0.234
IN 3	0.1234
IN 4	0.654
IN 5	0.789
IN 6	1.345
IN 7	0.8394
IN 8	1.96543
IN 9	2.4565
IN 10	0.5747

OUT

1.3437962907720776E35

ATTIVA

The interface



Regression function

Neural Network - V 3.0 - 1.gnn

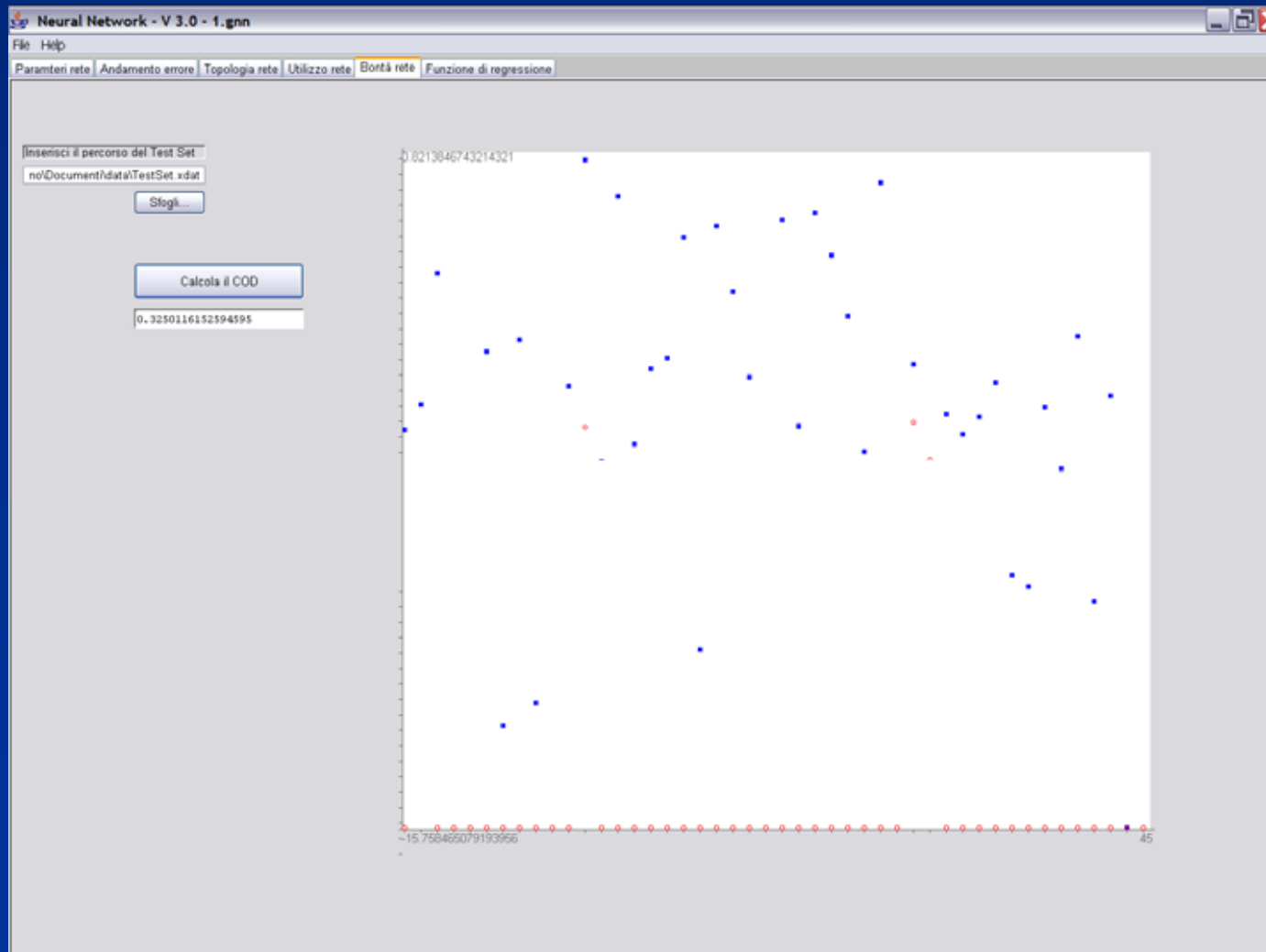
File Help

Parametri rete Andamento errore Topologia rete Utilizzo rete Bontà rete Funzione di regressione

Calcola la funzione di regressione

```
RegFuncLay_6 = [603.5952661749288]n0,0 + [130.12413431394362]n0,6 + [[-1409.491053446148]]c +  
[3477.6799312795612]n0,2 + [-1506.2879361468772]n0,3 + [0.11144138373248324]n1,0 +  
[-1.9538273589384272]n1,1 + [39.643341321644584]n0,4 + [1.142277371540958]n1,2 +  
[-0.029826657770687604]n2,0 + [-0.13246811018507864]n2,4 + [-3.746967856548508]n0,8 +  
[0.018816629709426540]n1,4 + [-201.89600986145024]n0,5 + [0.72170611125177776]n1,3 +  
[0.16352768024321306]n2,2 + [-0.21157674156507008]n2,5 + [0.16728653257169412]n3,0 +  
[-0.27482023860500116]n3,2 + [0.10041322204386382]n2,1 + [19.055780325726344]n0,9 +  
[0.082876557060385984]n1,5 + [0.63158365180508856]n2,3 + [0.09840791197610522]n3,1 +  
[-0.10167686654757358]n3,3 + [-2.003075842783616]n4,0 + [0.8569968442006512]n4,1 +  
[-1.6055583542767536]n3,4 + [-0.067003603344502064]n2,6 + [0.8967898259329656]n3,5 +  
[-0.88971302987761408]n4,2 + [0.48146785431202056]n4,3 + [-0.9083268984480136]n5,0 +  
[0.46764878352640872]n5,1
```

Computing evaluation parameters



Model Validation

- Total variance:
- Explained variance:
- Not explained variance:
- Determination coefficient COD
- COD is $0 \leq r^2 \leq 1$:

$$\frac{\sum_{i=1}^n (y_i - m_y)^2}{n}$$
$$\frac{\sum_{i=1}^n (\hat{y}_i - m_y)^2}{n}$$
$$\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$
$$\frac{\sum_{i=1}^n (\hat{y}_i - m_y)^2}{\sum_{i=1}^n (y_i - m_y)^2}$$

Standard

- *XML (Extensible Markup Language).*
- Standard interchange language
- Trained net saved in XML
- Net parsed in the application
- Open to future web services and applications

dataset XML

- `<?xml version="1.0" encoding="UTF-8" standalone="yes" ?>`
- `<?xml-stylesheet type="text/xsl" href="nostrisenz1rigOut.xls.xsl"?>`
- `<table name="Foglio1$">`
- `<Foglio1>`
- `<MW>158.27</MW>`
- `<nAT>29</nAT>`
- `<nBT>28</nBT>`
- `<MR>44.877</MR>`
- `<PSA>17.07</PSA>`
- `<GMTIV>1509</GMTIV>`
- `<SPI>25.632</SPI>`
- `<piID>76</piID>`
- `<SRW08>700</SRW08>`
- `<SRW10>2520</SRW10>`
- `<Output>0.574966829</Output>`
- ...
- `</table>`

Root
element

name for all
the children

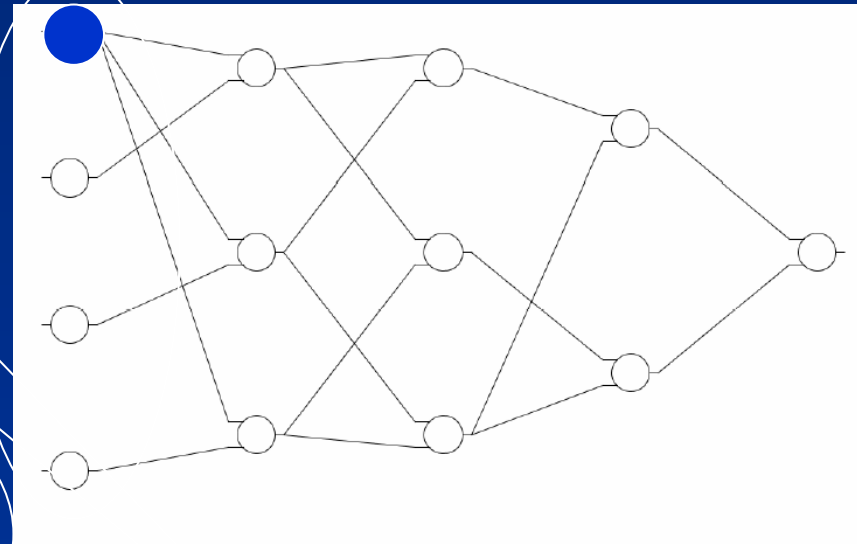
Tuple of
dataset

Features

Network in XML

root Layer Neuron

- `<?xml version="1.0" encoding="UTF-8" ?>`
- `<NeuralNetwork>`
- `<Layer id="0">`
- `<Neuron id="0">`
- `<CombinationFunction name="FirstLayerNeuron" />`
- `<ActivationFunction name="IvakhnenkoPolynomial" num-params="6">`
- `<AFPParameter id="1">0.0</AFPParameter>`
- `<AFPParameter id="2">0.0</AFPParameter>`
- `<AFPParameter id="3">0.0</AFPParameter>`
- `<AFPParameter id="4">1.0</AFPParameter>`
- `<AFPParameter id="5">0.0</AFPParameter>`
- `<AFPParameter id="6">0.0</AFPParameter>`
- `</ActivationFunction>`
- `</Neuron>`



Sub-element of Neuron

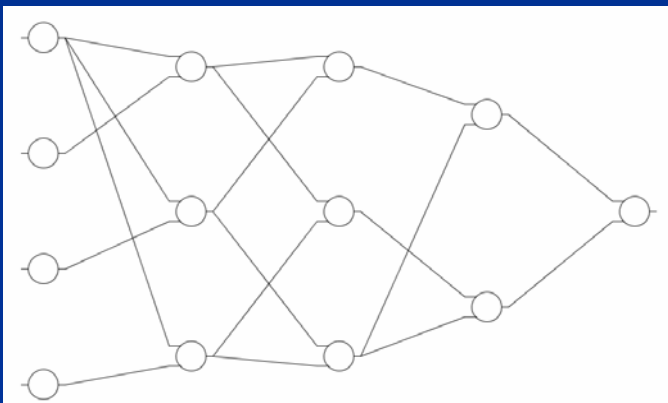
Parameters of the activation function

Data, equations, models as text (XML)

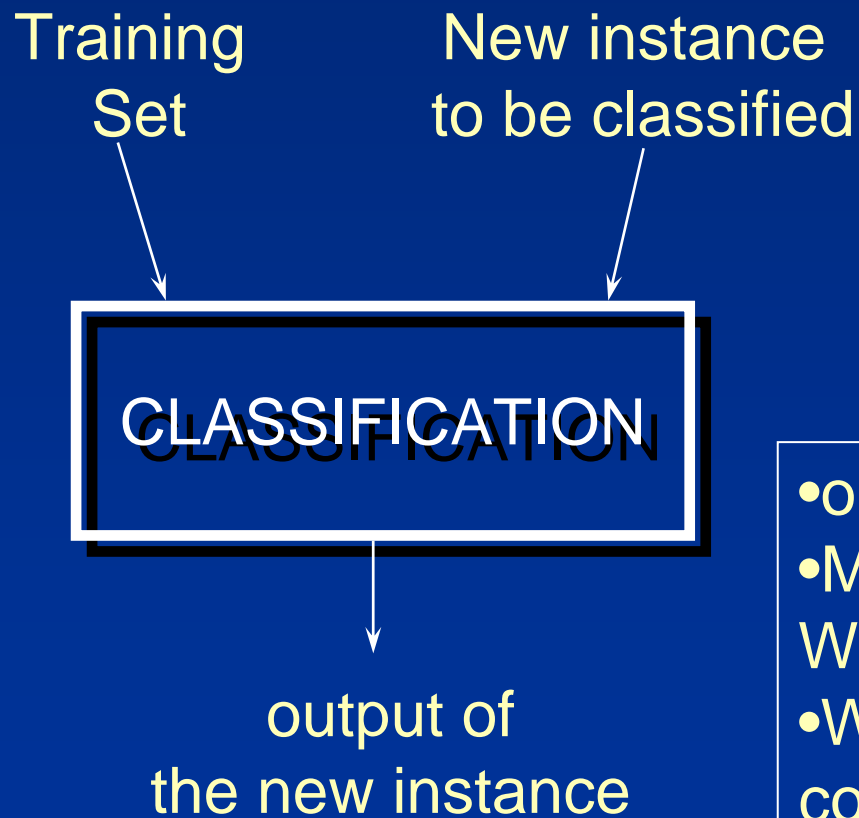
Emphasis on
services

```
<?xml version="1.0" encoding="UTF-8" ?>
<NeuralNetwork>
<Layer id="0">
<Neuron id="0">
<CombinationFunction name="FirstLayerNeuron" />
<ActivationFunction name="IvakhnenkoPolynomial" num-params="6">
<AFPParameter id="1">0.0</AFPParameter>
<AFPParameter id="2">0.0</AFPParameter>
<AFPParameter id="3">0.0</AFPParameter>
<AFPParameter id="4">1.0</AFPParameter>
<AFPParameter id="5">0.0</AFPParameter>
<AFPParameter id="6">0.0</AFPParameter>
</ActivationFunction>
</Neuron>
```

```
= "2">
tionFunction name="NeuronSelection" num-params="2">
eter id="0">
" type="Nref" lr="0" nrm="0" />
" type="weight">1</Arg>
meter>
eter id="1">
" type="Nref" lr="0" nrm="3" />
" type="weight">1</Arg>
meter>
ationFunction>
hFunction name="IvakhnenkoPolynomial" num-params="6">
eter id="1">0.0895283684330046</AFPParameter>
eter id="2">0.27905286981012445</AFPParameter>
eter id="3">0.7434228721524345</AFPParameter>
<AFPParameter id="4">2.610780910353352</AFPParameter>
<AFPParameter id="5">-1.730970145378251</AFPParameter>
<AFPParameter id="6">-0.35500543418743646</AFPParameter>
</ActivationFunction>
</Neuron>
```



Next step



- one model is built
- More models are similar. Which is the best?
- What if (a few) models are combined?

What is ensemble learning?

Ensemble learning refers to a collection of methods that learn a target function by training a number of individual learners and combining their predictions

Hybrid – a technical definition

- involves the use of two or more intelligent techniques and approaches, (neural networks, knowledge-based methods, fuzzy techniques, genetic algorithms, agent-based techniques, case based reasoning etc).
- The combination can be done in any form, either by a modular integration of two or more intelligent methodologies, or by fusing one methodology into another...

Ensembles: different names

- multiple models
 - multiple classifier systems
 - combining classifiers (regressors *etc*)
 - integration of classifiers
 - mixture of experts
 - decision committee
 - committee of experts
 - classifier fusion
 - ...
- base classifiers
 - component classifiers
 - individual classifiers
 - members (of a decision committee)
 - level 0 experts
 - ... what else?

ensembles

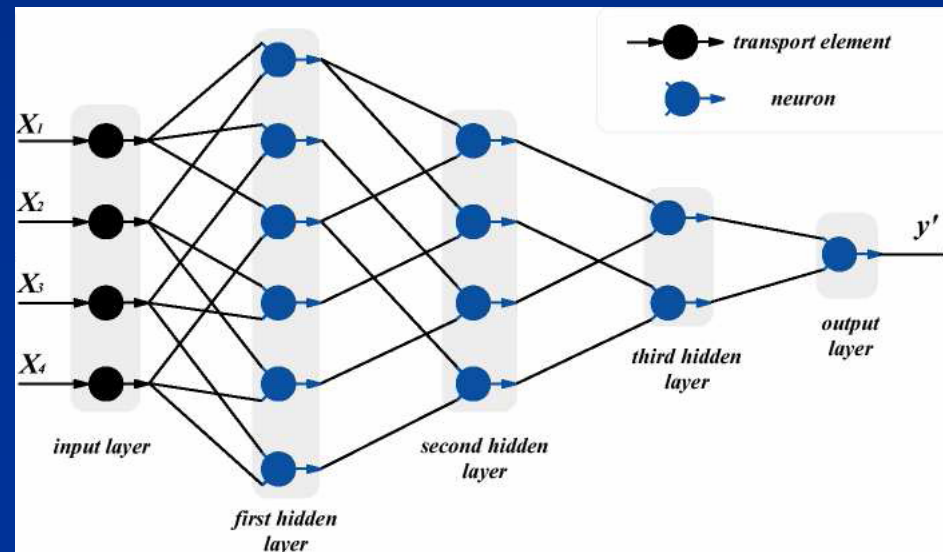


classifiers in ensembles

2 views on ensembling

- A method to reduce errors (bias and variance)
- A method to integrate more concepts (hybrid)

if a then b



Why ensemble learning?

- *Accuracy*: a more reliable result by combining the output of multiple "experts"
- *Efficiency*: a complex problem can be decomposed into multiple sub-problems that are easier to understand and solve
- no single model that works for all pattern recognition problems
- "To solve really hard problems, we'll have to use several different representations..... It is time to stop arguing over which type of pattern-classification technique is best..... Instead we should work at a higher level of organization and discover how to build managerial systems to exploit the different virtues and evade the different limitations of each of these ways of comparing things." Minsky, 1991.

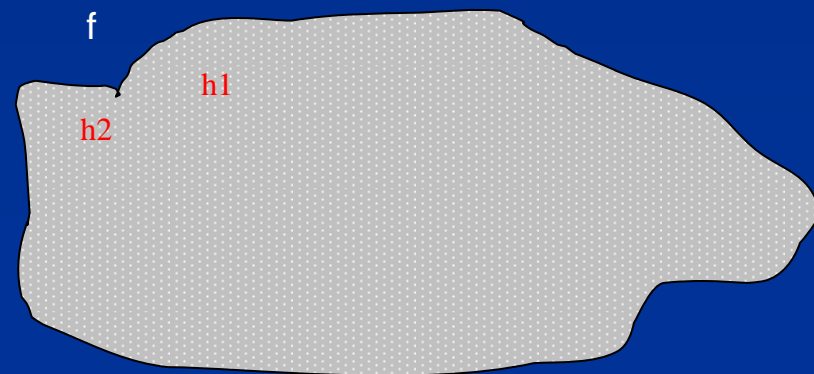
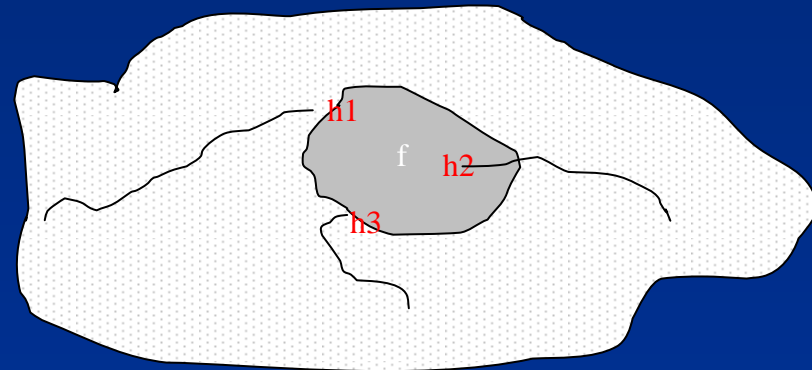
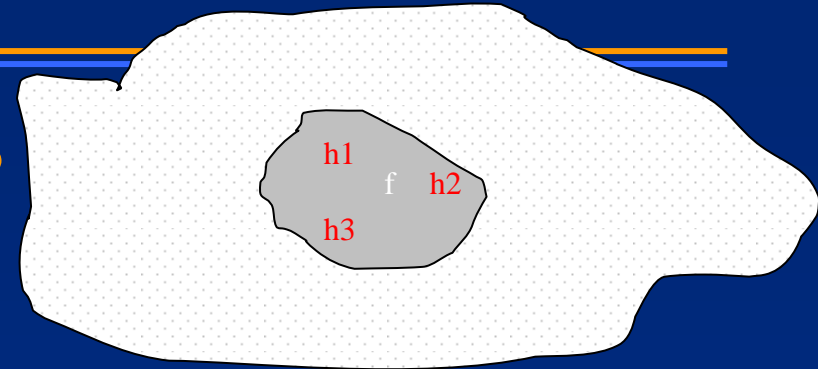
When ensemble learning?

- When you can build base classifiers that are *more accurate than chance*, and
- that are as much as possible *independent* from each other

Why do ensembles work?

ensembles overcome three problems:

- **Statistical Problem** : there are many hypotheses with the same accuracy, the learning algorithm chooses one of them – better to mix them.
- **Computational Problem**: the learning algorithm cannot guarantee reaching the best hypothesis.
- **Representational Problem**: the hypothesis space does not contain any good approximation of the target classes.



Theoretical results

Hansen & Solomon (1990):

If we can assume that classifiers are independent in predictions and their accuracy $> 50\%$, can push accuracy arbitrarily high by combining more classifiers

Key assumption:

classifiers are *independent* in their predictions

How to make an effective ensemble?

Two basic decisions when designing ensembles:

1. How to generate base classifiers?
(*generation strategy*)
2. How to integrate them?
(*integration strategy*)

Methods for Independently Constructing Ensembles

One way to force a learning algorithm to construct multiple hypotheses is to run the algorithm **several times** and provide it with somewhat **different data in each run**.

This idea is used in methods as:

- *Bagging (instance selection)*
- *Feature-Selection Ensembles*
- *Random Forest* - both sources of randomness

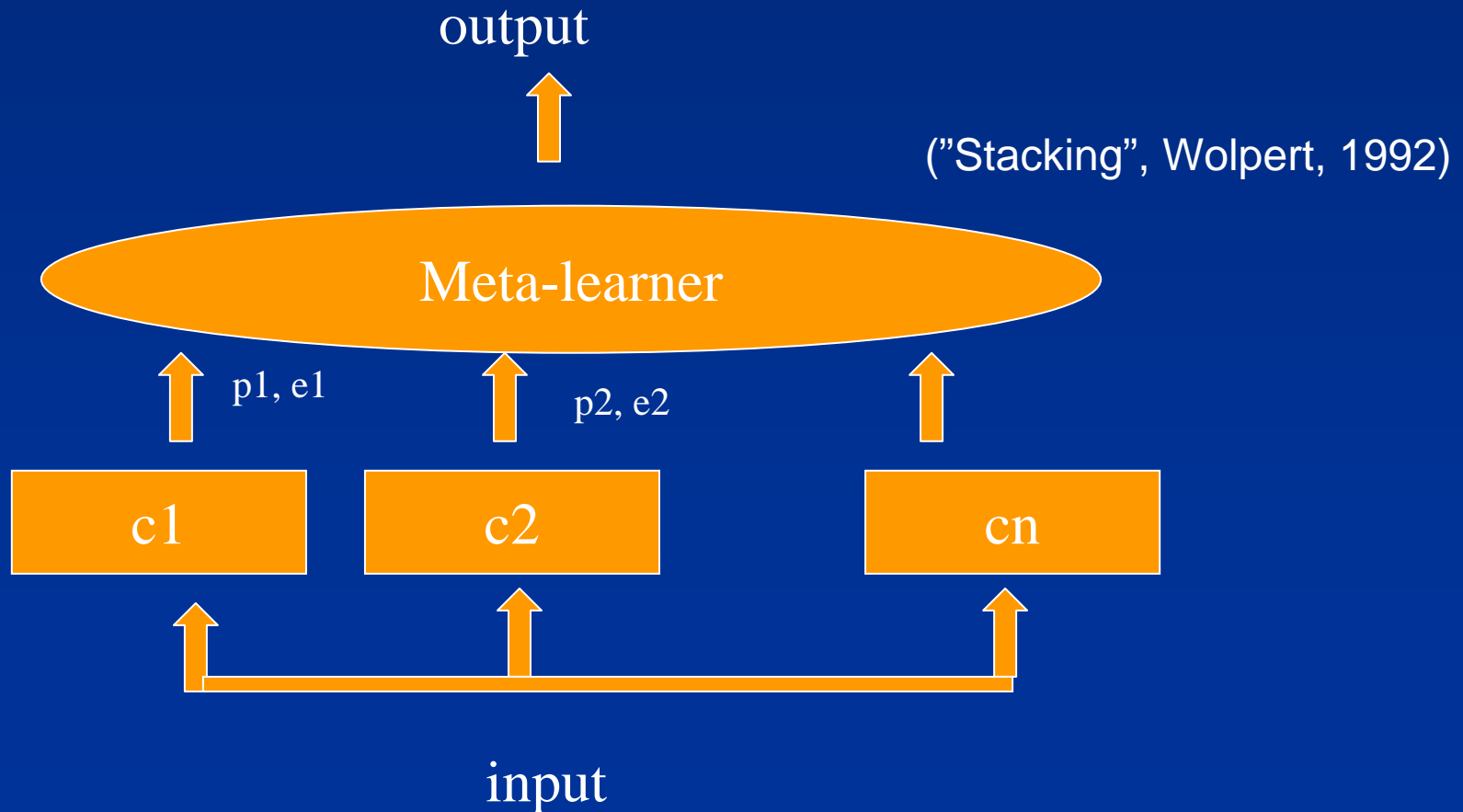
Methods for Coordinated Construction of Ensembles

The key idea is to learn complementary classifiers so that instance classification is realized by taking **weighted sum** of the classifiers. This idea is used in :

- Boosting - new models are influenced by performance of previously built ones
- Stacking.

Stacking

Uses *meta learner* instead of voting
Predictions of base learners (*level-0 models*) are used as
input for meta learner (*level-1 model*)



Overfitting in ensembles

- Not that much research has been done to this time
- A surprising finding:
ensembles of overfit base classifiers
(DTs, ANNs) are in many cases better than the
ensembles of non-overfit base classifiers
 - This is related most probably to the fact that in that case the
ensemble diversity is much higher

acknowledgements

- The development of poliGMDH is gratefully acknowledged to different EU projects:
 - Allelochem (2003-5) – the power of GMDH has been proven on a QSAR problem
 - EASYRING (2003-6) – a complete working version of poliGMDH (Neural Network v3) and the development of models for endocrine disruptors
 - Demetra (2003-6)– the power of the method and its use in ensembling
 - ION (2005-6) – the role of the method in the process of drug discovery