

# STRUCTURAL OPTIMIZATION OF REGRESSION MODELS BASED ON A GENETIC ALGORITHMS



Speakers:

**Roman Kvyetnyy**, DrSc in Engineering, Professor, Corresponding Member of NAES of Ukraine, Professor of the Department of Automation and Intellectual Information Technologies

**Yaroslav Ivanchuk**, DrSc in Engineering, Professor, Professor of the Department of Computer Science

**Oleksii Kozlovskiy**, Postgraduate student, Department of Computer Science

# PROBLEM STATEMENT

$$\hat{y} = f(\mathbf{x}; \boldsymbol{\theta}, \mathbf{s})$$

$D = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  - training sample;

$\mathbf{x}_i \in \mathbb{R}^p$  - vector of input features;

$y_i \in \mathbb{R}$  - observed response value;

$f \in F_s$  - class of admissible regression functions;

$\boldsymbol{\theta} \in \Theta(\mathbf{s})$  - vector of model tuning parameters;

$\mathbf{s} \in S$  - vector of structural variables that determine the structure of regression model;

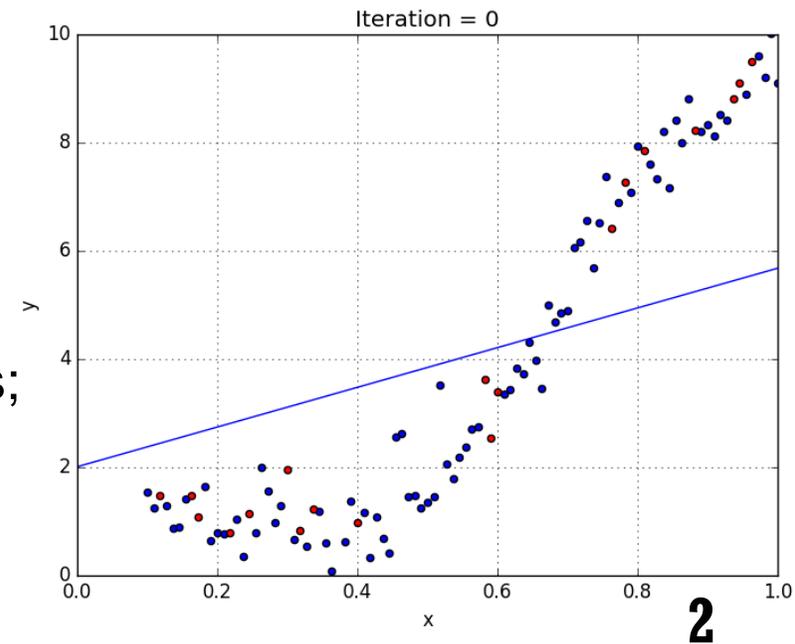
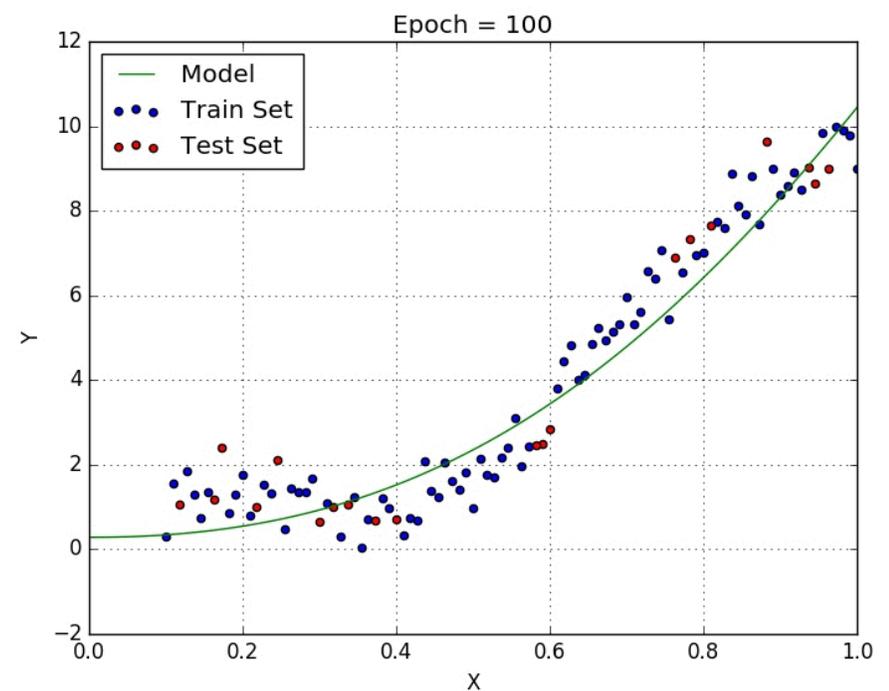
## Regression model training task:

$$\frac{1}{N} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i; \boldsymbol{\theta}, \mathbf{s})) + \lambda \Omega(\boldsymbol{\theta}, \mathbf{s}) \rightarrow \min_{\mathbf{s} \in S} \min_{\boldsymbol{\theta} \in \Theta}$$

$L(y_i, f(\mathbf{x}_i; \boldsymbol{\theta}, \mathbf{s})) = (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}, \mathbf{s}))^2$  - mean squared error (MSE) loss;

$\Omega(\boldsymbol{\theta}, \mathbf{s})$  - regulatory functionality;

$\lambda$  - regularisation coefficient.



# AIM AND OBJECTIVES

- *The aim of this work* is to improve the accuracy of predicting target variable values based on observed data by developing a structural optimization method using evolutionary computation.

## Research tasks:

- to develop a genetic algorithm-based optimization method that determines the optimal architecture and type of regression model in high-dimensional feature spaces by sequentially selecting, combining, and varying candidate solution parameters using evolutionary mechanisms.
- to develop an approach for identifying the operational parameters of a genetic algorithm using a benchmark nonlinear function, which will enable the effective detection of the optimal structure of a regression model;
- to assess the effectiveness of established regression model types by applying the developed genetic algorithm-based structural optimization method to benchmark dataset models that represent the descriptive characteristics of real-world systems.

# GENETIC ALGORITHM-BASED STRUCTURAL OPTIMIZATION METHOD FOR REGRESSION MODELS

$$\mathbf{c} = (s, \boldsymbol{\theta}) \quad \begin{array}{l} S - \text{chromosomes encoding the model structure;} \\ \boldsymbol{\theta} - \text{genes directly encoding the regression model parameters.} \end{array} \quad (1)$$

The chromosome value as a solution for searching the optimal number of the most important features:

$$\mathbf{c} \equiv \mathbf{s} = (s_1 \ s_2 \ \dots \ s_p) \in B^p, \ B = \{0, 1\} \quad (2)$$
$$s_j = \begin{cases} 1, & \text{if the } j\text{-th feature is selected;} \\ 0, & \text{if the } j\text{-th feature is not selected,} \end{cases}$$

Fitness function: 
$$F(\mathbf{c}) = -J(\mathbf{c}) = -\left[ \frac{1}{N} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i^{(c)}; \boldsymbol{\theta}^*(\mathbf{c}))) + \lambda \sum_{j=1}^p c_j \right]. \quad (3)$$

$\sum_{j=1}^p c_j$  – the number of selected features;  
 $\lambda$  – the penalty value for an excessive number of features.

Structural optimization of a regression model is the search for a solution to another fitness function optimization problem:

$$\mathbf{c}^* := \arg \max_{\mathbf{c} \in \{0,1\}^p} F(\mathbf{c}) \quad (4)$$

Population: 
$$\mathbf{P}^{(m)} = (\mathbf{c}_1^{(m)}, \mathbf{c}_2^{(m)}, \dots, \mathbf{c}_M^{(m)}) \in \mathbf{C}^M = (\{0,1\}^p)^M \quad (5)$$

$m$  – generation number,  $M$  – population size.

# GENETIC ALGORITHM FOR STRUCTURAL OPTIMIZATION OF A REGRESSION MODEL

$$P^{(m+1)} = E(P^{(m)}) \cup M(\aleph(S(P^{(m)}))) \quad (6)$$

$$E(P^{(m)}) = \arg \max_{\varepsilon} \{F(\mathbf{c}) \mid \mathbf{c} \in P^{(m)}\} \quad (7)$$

$\varepsilon \in \{1, 2, \dots, M\}$  - number of elite individuals;

►  $P_{sl}^{(m)} = S(P^{(m)}) = \{\mathbf{c}_1^{sl}, \mathbf{c}_2^{sl}, \dots, \mathbf{c}_M^{sl}\}$  - reproduction (selection) operator; (8)

$$\mathbf{c}_j^{sl} = \arg \max_{\mathbf{c} \in T_j} F(\mathbf{c}) \text{ - independent tournament of size } \tau;$$

►  $P_{cr}^{(m)} = \aleph(P_{sl}^{(m)}) = \bigcup_{j=1}^{M/2} \aleph(\mathbf{c}_{2\delta-1}, \mathbf{c}_{2\delta}; r_{\delta})$  - genetic crossover operator; (9)

$r_{\delta} \in \{1, 2, \dots, p-1\}$  - randomly selected crossover point for the  $\delta$ -th pair of parent chromosomes;

$\aleph(\mathbf{c}_{2\delta-1}, \mathbf{c}_{2\delta}; r_{\delta}) = \mathbf{c}'_{2\delta-1}, \mathbf{c}'_{2\delta}$  - two offspring chromosomes are generated according to a given rule.

►  $P_{mut}^{(m)} = M(P_{sl}^{(m)}) = \{M(\mathbf{c}_j) \mid \mathbf{c}_j \in P^{(m)}\}$  - genetic mutation operator; (10)

$$M(\mathbf{c}_j) = \mathbf{c}'_j = (\mathbf{c}'_{1j}, \mathbf{c}'_{2j}, \dots, \mathbf{c}'_{ij}) \quad \mathbf{c}'_{ij} = \begin{cases} 1 - c_{ij}, & \text{with probability } p_{mut}; \\ c_{ij}, & \text{with probability } 1 - p_{mut}. \end{cases}$$

# REGRESSION METHODS

## LINEAR REGRESSION

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

$\mathbf{w}$ ,  $b$  – model parameter vectors

## k-NEAREST NEIGHBORS (kNN)

$$f(\mathbf{x}) = \frac{1}{k} \sum_{j \in N_k(\mathbf{x})} y_j$$

$k$  – number of nearest points,

$N_k(\mathbf{x})$  –  $k$  nearest points

## GRADIENT BOOSTING REGRESSION (GBR)

$$f(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x})$$

$\alpha_t(\mathbf{x})$  – weight coefficient of the  $t$ -th decision tree

## RANDOM FOREST (RF)

$$f(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T h_t(\mathbf{x})$$

$T$  - number of decision trees,

$h_t(\mathbf{x})$  - predicted value from the  $t$ -th decision tree

## MULTILAYER PERCEPTRON (MLP)

$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

$L$  - Number of network layers

$$y_i = 10 \sin(\pi x_{i1} x_{i2}) + 20(x_{i3} - 0,5)^2 + 10x_{i4} + 5x_{i5} + \varepsilon_i; \quad \mathbf{X} = \begin{bmatrix} x_{ij} \end{bmatrix} \in \mathbb{R}^{N \times p} \quad (11)$$

$x_{i1}, \dots, x_{i5}$  - relevant features;  $x_{i6}, \dots, x_{ip}$  - noisy features;  $\varepsilon_i$  - additive noise;

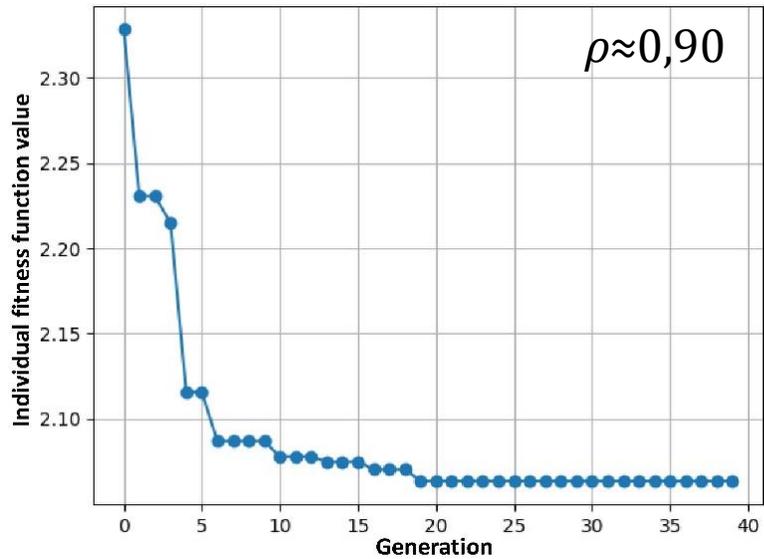
**Table 1 - The target variable value based on relevant and noisy features**

No	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_{17}$	$x_{18}$	$x_{19}$	$x_{10}$	$x_{11}$	$x_{12}$	$x_{13}$	$x_{14}$	$x_{15}$	$y$
0	0,542	0,601	0,558	0,735	0,953	0,854	0,623	0,875	0,320	0,275	0,032	0,242	0,845	0,119	0,774	20,72
1	0,537	0,377	0,939	0,269	0,289	0,597	0,638	0,047	0,152	0,904	0,602	0,677	0,083	0,489	0,697	13,93
2	0,217	0,488	0,488	0,068	0,565	0,702	0,546	0,534	0,401	0,312	0,971	0,608	0,653	0,592	0,048	6,770

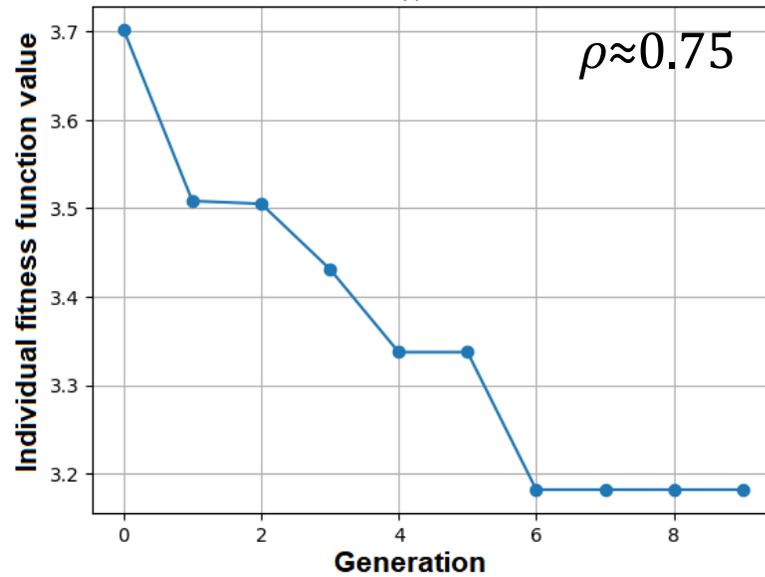
**Table 2 - Genetic algorithm parameters for the Friedman-1 dataset**

Type of regression model	Population size, $M$	Crossover/mutation probability, $p_{cr}/p_{mut}$	Penalty value, $\lambda$
kNN	14-50	0,7/0,2	0,01
RF	16-60	0,9/0,1	0,01
GBR	18-40	0,8/0,1	0,04
LR	16-40	0,8/0,2	0,02
MLP	18-40	0,9/0,2	0,01

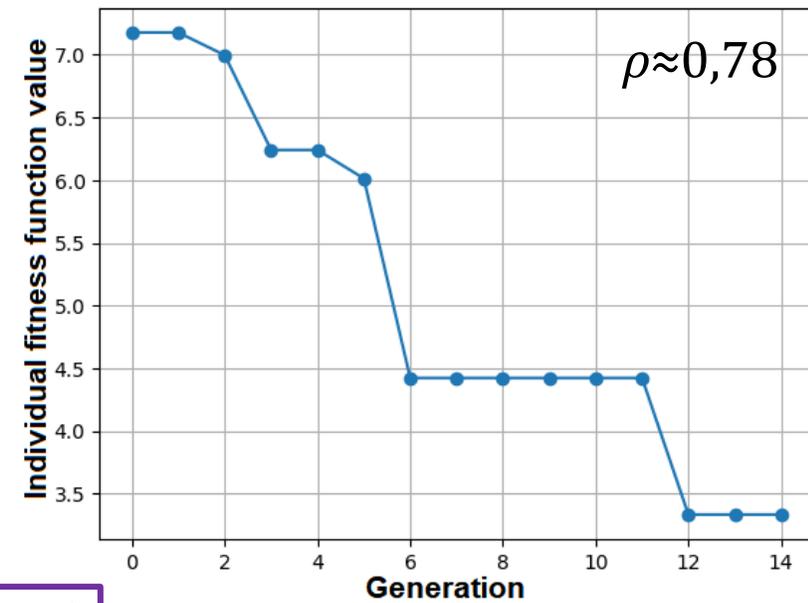
# Gradient Boosting Regression



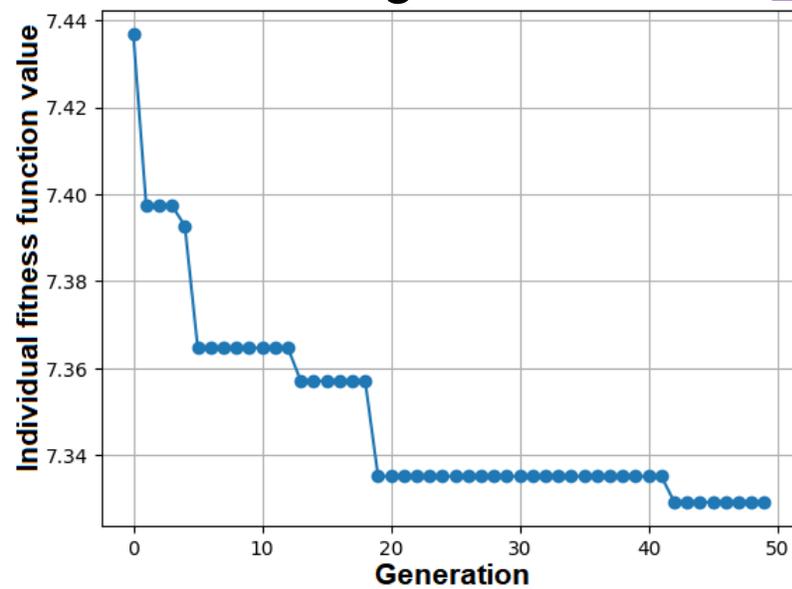
# Random Forest



# kNN

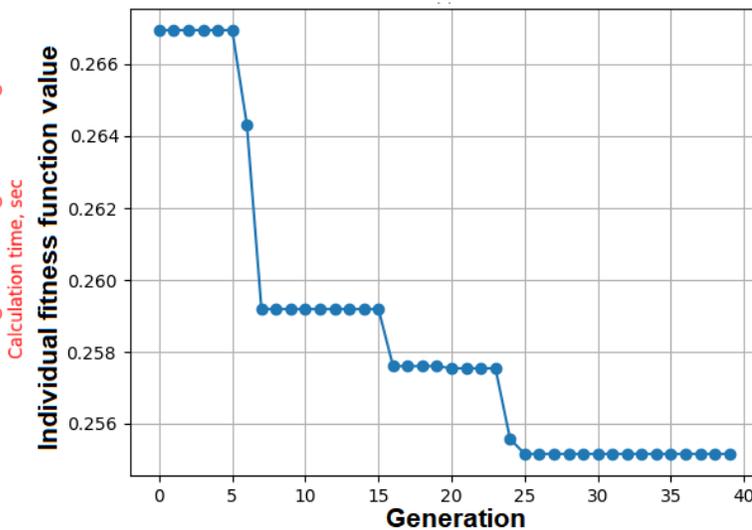


# Linear regression

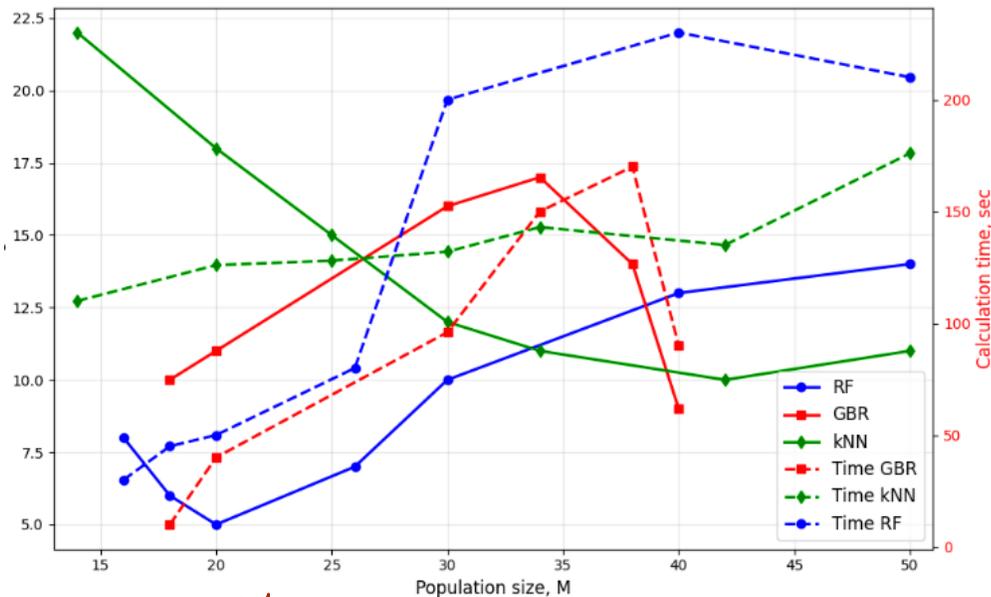


$$s = (x_1 \ x_2 \ \dots \ x_{15}) = (1111100000000000)$$

# MLP



$$(111110000100100)$$



$$(111111000100000)$$

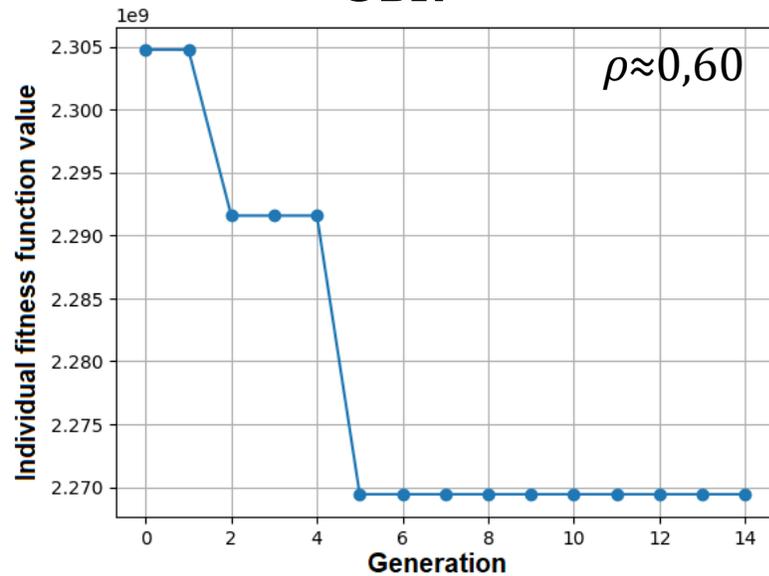
# EXPERIMENTAL RESULTS

Table 3. Response variable of the California Housing dataset  $\rho \approx 0,90$

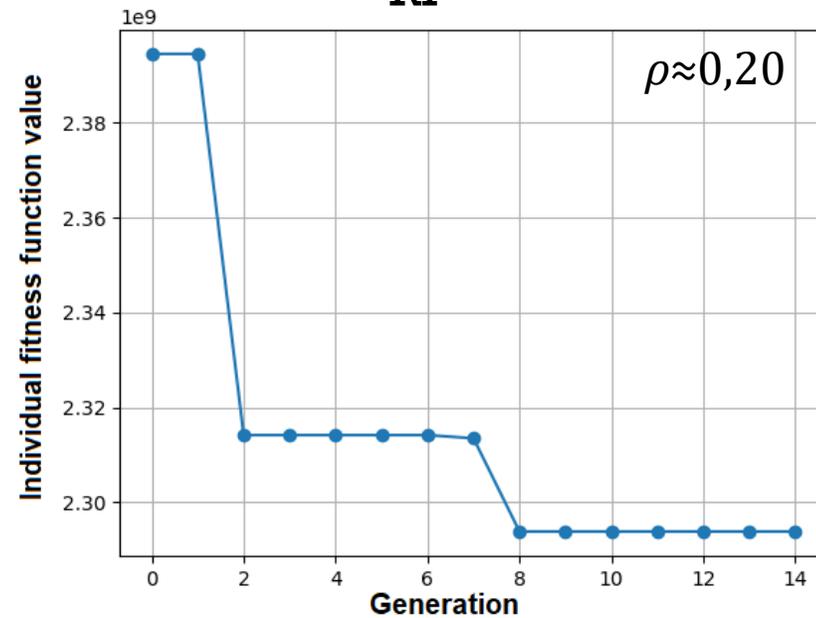
No	Longitude	Latitude	Housing_median_age	Total_rooms	Total_bedrooms	Population	Households	Median_income	Median_house_value
	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$	$X_7$	$X_8$	Y
0	-114,31	34,19	15,0	5612,0	1283,0	1015,0	472,0	1,4936	66900,0
1	-114,47	34,40	19,0	7650,0	1901,0	1129,0	463,0	1,8200	80100,0
2	-114,56	33,69	17,0	720,0	174,0	33,0	117,0	1,6509	85700,0

$$s = (x_1 \ x_2 \ \dots \ x_8) \equiv (s_1 \ s_2 \ \dots \ s_8)$$

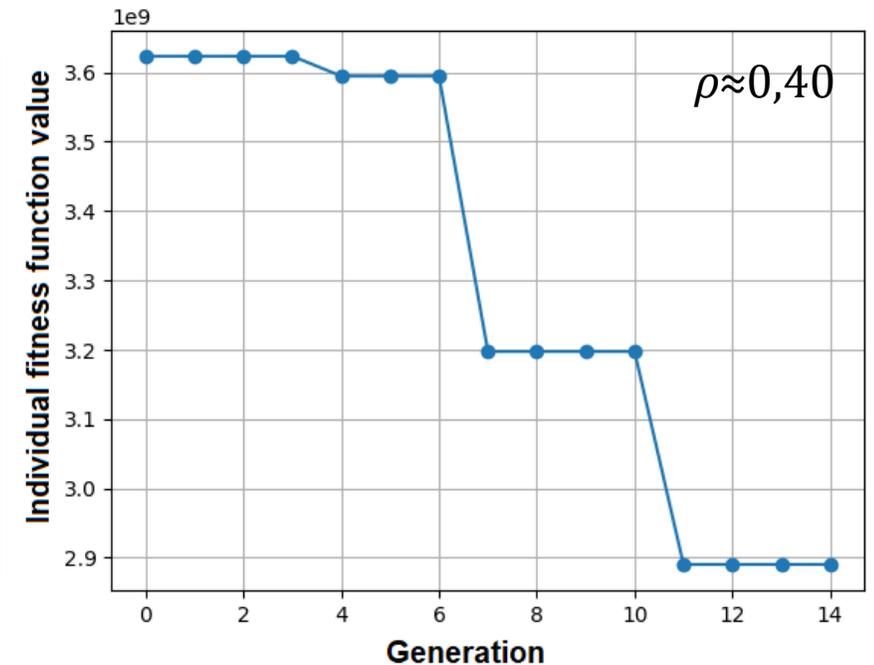
### GBR



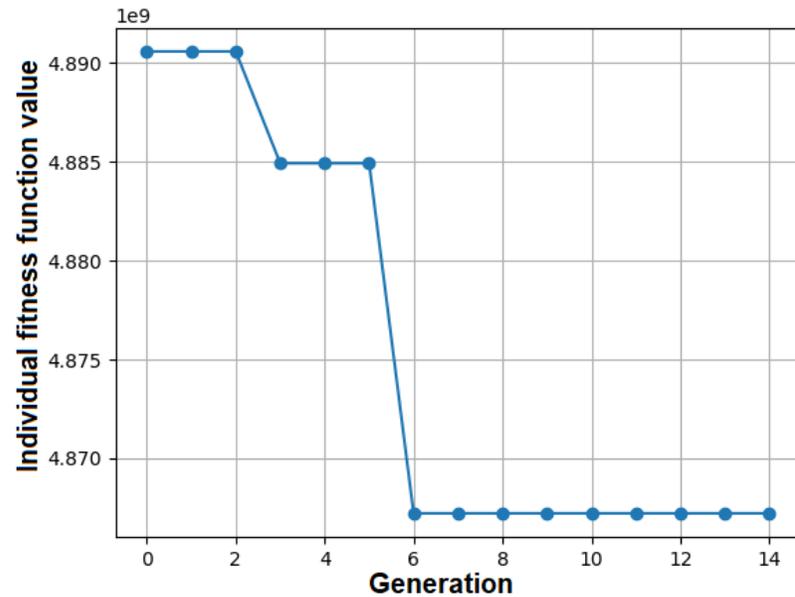
### RF



### kNN



LR



MLP

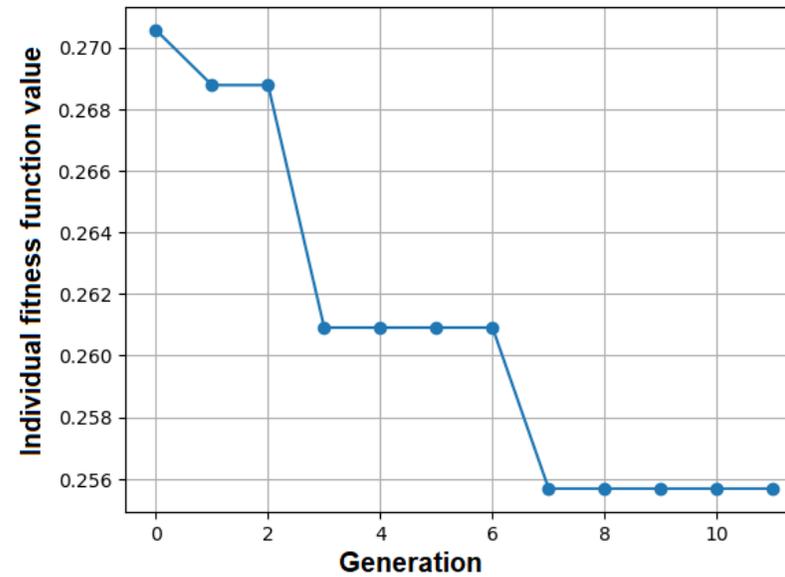
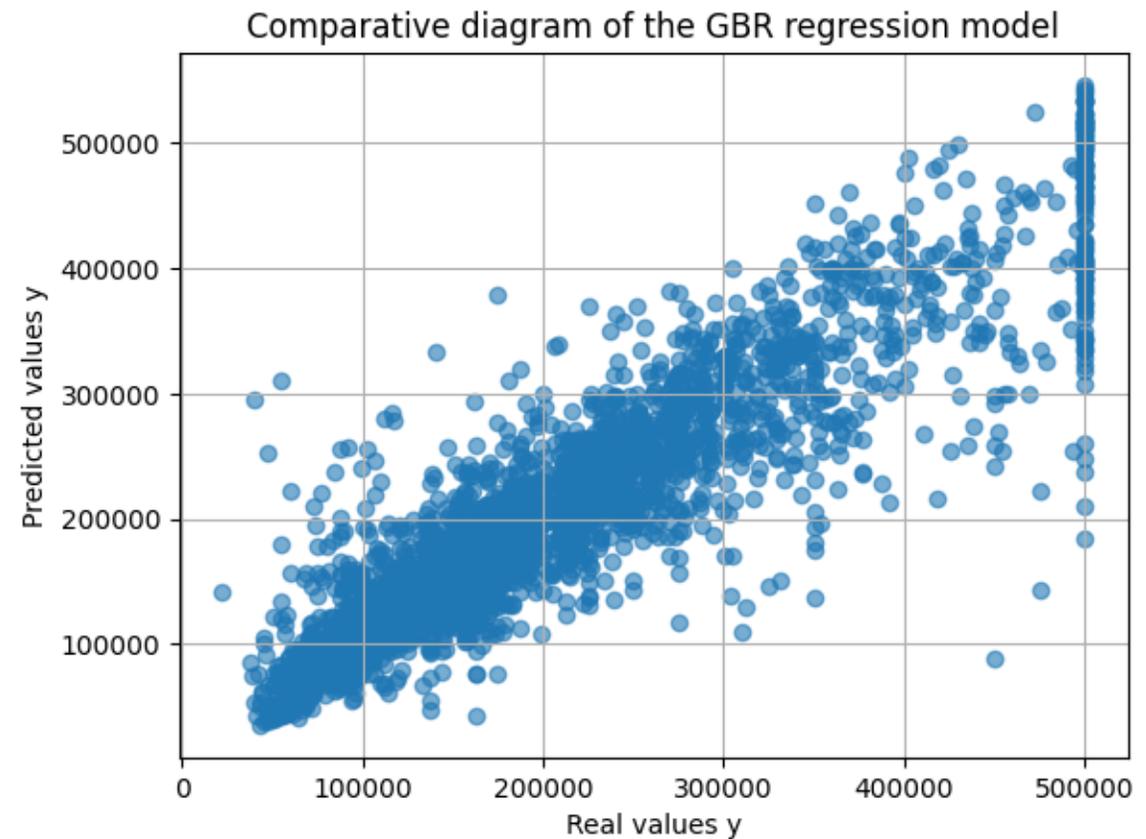


Table 4. Results of regression model structural optimization using a genetic algorithm

Regression Model	Vector of optimal features			Coefficient of determination of the baseline regression model, $R^2$	Coefficient of determination after structural optimization, $R^2$		
	Genetic Algorithm	Lasso	Recursive Feature Elimination (RFE)		Genetic Algorithm	Lasso	Recursive Feature Elimination (RFE)
kNN	[1,1,0,0,0,0,0,0]	[1,1,1,1,1,1,1,1]	[1,1,0,0,1,1,0,1]	0,70	0,77	0,70	0,70
RF	[1,1,0,1,0,1,0,1]			0,80	0,82	0,80	0,81
GBR	[1,1,1,1,0,1,0,1]			0,81	0,83	0,81	0,81
LR	[1,1,1,1,1,1,1,1]			0,62	0,62	0,62	0,60
MLP	[1,1,1,1,1,1,1,1]			0,75	0,75	0,75	0,74

# REFERENCES

- The paper addresses the solution of an important scientific and practical problem of structural optimization of regression models (SORM) using a genetic algorithm (GA). The application of the GA increased the prediction accuracy of regression models by an average of 5% compared to the use of regularization methods.
- A conditional multi-objective SORM problem is formulated using a two-level minimization of the empirical risk functional on a training dataset for ensemble, local, and globally parametric regression methods. The problem is solved using a genetic algorithm with a fitness function based on a quadratic loss and a feature-number penalty.
- A parameter identification approach for the genetic algorithm based on the nonlinear Friedman-1 function was developed to determine the optimal population size, crossover and mutation probabilities, and the fitness penalty value.
- The testing results demonstrated an average increase in the coefficient of determination of the regression models by  $\Delta R^2 = 0.036$ , while the number of irrelevant features was reduced by an average of 46% for ensemble and local regression methods based on the genetic algorithm. In addition, the results of the GA iterative computations in SORM exhibited linear-type asymptotic convergence with an empirical coefficient of  $\rho \approx 0.90$ .



# Thank you for your attention!



Vinnytsia National  
Technical University



***Roman Kvyetnyy,***  
DrSc in Engineering,  
Professor,  
Corresponding Member  
of NAES of Ukraine,  
Professor of the  
Department of  
Automation and  
Intellectual Information  
Technologies



***Yaroslav Ivanchuk,***  
DrSc in Engineering,  
Professor, Professor of  
the Department of  
Computer Science



***Oleksii Kozlovskyy,***  
Postgraduate student,  
Department of Computer  
Science