

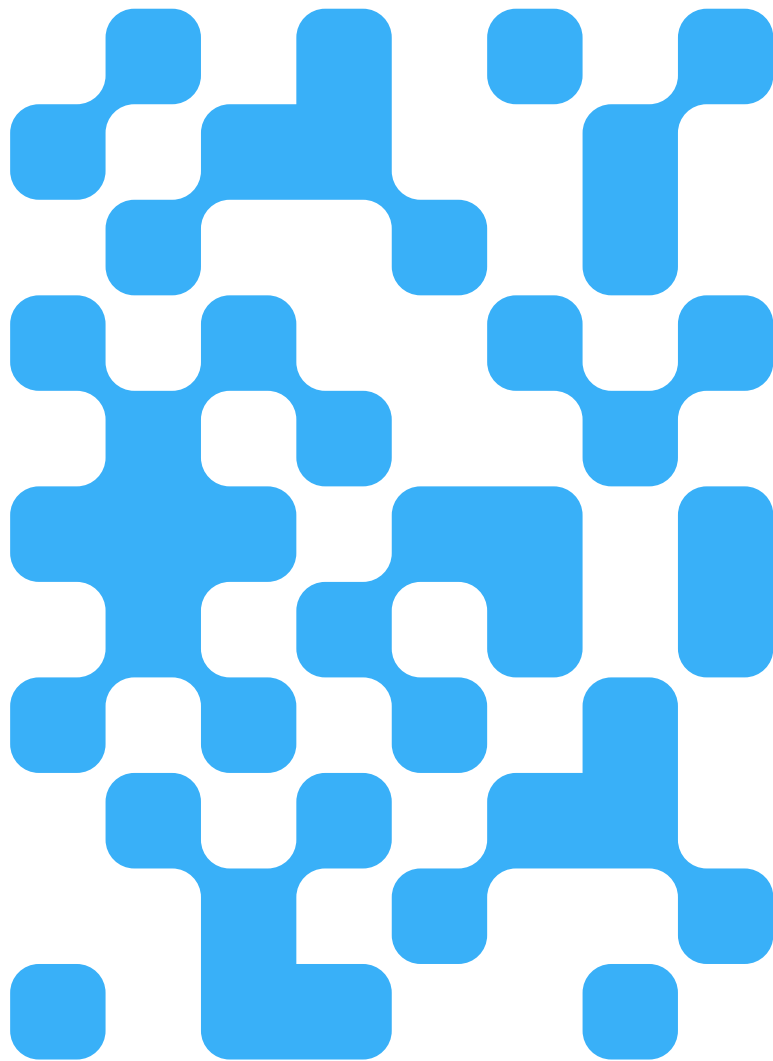


# Machine Learning

2024 (ML-2024)

Lecture 2. Linear models

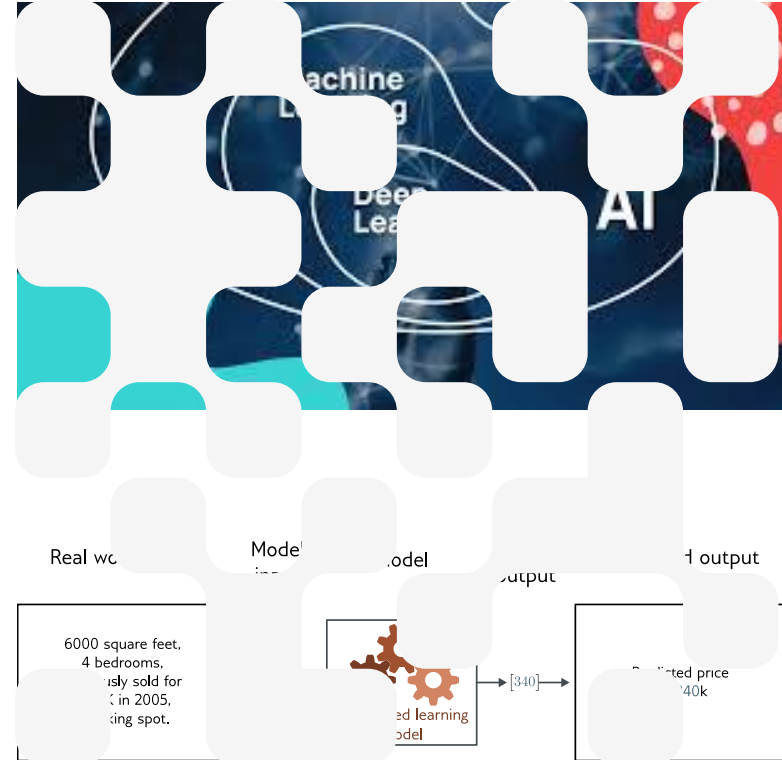
by Alexei Valerievich Kornaev, Dr. habil. in Eng. Sc.,  
Researcher at the RC for AI, Assoc. Prof. of the Robotics and CV  
Master's Program, [Innopolis University](#)  
Researcher at the RC for AI, [National RC for Oncology n.a. NN Blohin](#)  
Professor at the Dept. of Mechatronics, Mechanics, and Robotics,  
[Orel State University](#)



# Agenda

- I. Linear Regression and its Generalization
- II. Logistic Regression and its Generalization
- III. Setting of the models

All models are wrong, but some are useful.  
/George Box/



## Books

[Handbook on Machine Learning](#) by M. Artemyev et al.,  
Yandex, 2022 (in Russian)  
[Understanding Deep Learning](#) by Simon J.D. Prince, 2024  
[Practical Deep Learning / FastAI book](#) by Jeremy Howard  
[Deep Learning](#) by Ian Goodfellow and Yoshua Bengio and  
Aaron Courville, 2016.

## Online platforms, courses, resources

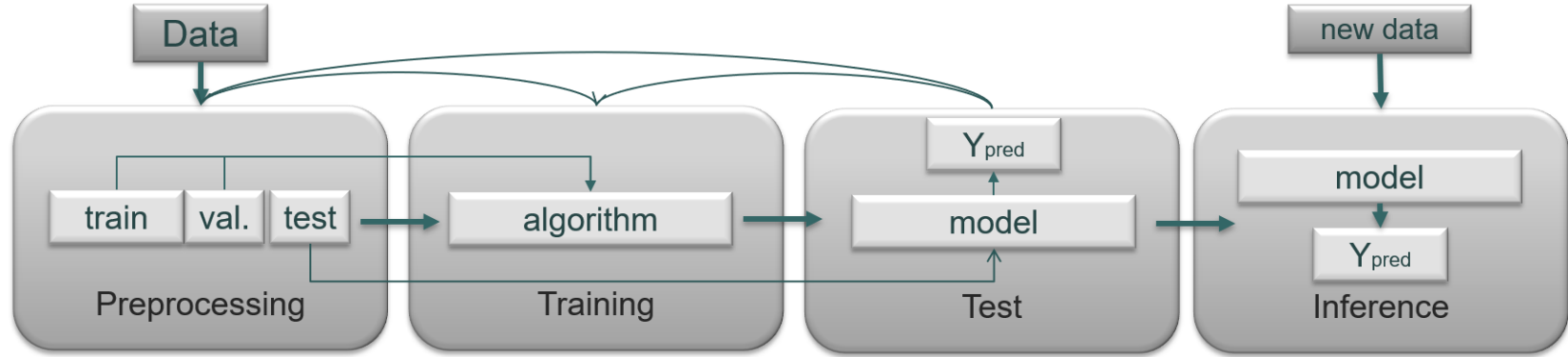
[Sirius](#) online courses on ML (in Russian)  
[Stepik](#) online courses (in Russian)  
[Hugging Face](#) online courses  
Coursera is unavailable so far

[MIT Introduction to Deep Learning](#), MIT, 2024  
[Lecture Hall of the Faculty of Applied Mathematics and Informatics](#) (in Russian)  
[Fast AI](#), courses, software, book by Jeremy Howard  
[Deep Learning](#), course by Semyon Kozlov (in Russian), 2019  
  
[3Blue1Brown](#), Animated Math  
[PyTorch Tutorial](#) by Patrick Loeber, 2020

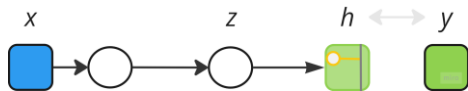
## #someLinks

Read here: <https://arxiv.org/>, <https://scholar.google.ru/>  
Collect the references here: <https://mendeley.com/>  
Draw here: <https://miro.com/app/dashboard/>  
Write the text here: <https://www.overleaf.com/project>  
Write the code here: <https://colab.research.google.com/>  
Collect the code here: <https://github.com/>  
Find the journal here: <https://journalfinder.elsevier.com/>  
Find the conference here: <https://portal.core.edu.au/conf-ranks/?search=A>

## Flowchart for an ML model design



# Linear Regression



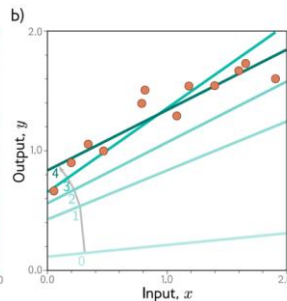
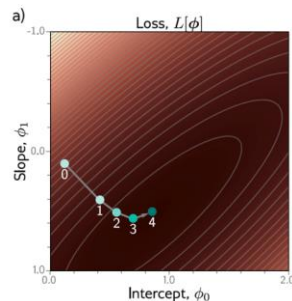
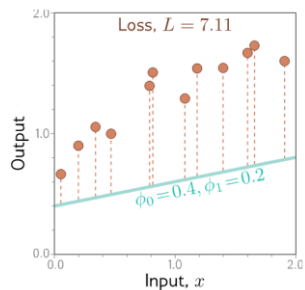
Model predicts output  $h$  given input  $x$

$$h = \phi_0 + \phi_1 x$$

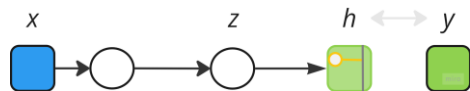
$$\mathbf{x} = \begin{bmatrix} x^{(1)} \\ \dots \\ x^{(m)} \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \dots \\ y^{(m)} \end{bmatrix}, \boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix}.$$

Consider a model  $f = [x^{(i)}, \boldsymbol{\phi}]$  parameterized with weights  $\boldsymbol{\phi}$  that maps each  $i$ -th input sample  $x^{(i)}$  into the output  $z^{(i)}$  which then transforms into the hypothesis  $h^{(i)}$  that should be close to the label  $y^{(i)}$ .

$$L(\boldsymbol{\phi}) = \frac{1}{2m} \sum_{i=1}^m (h^{(i)} - y^{(i)})^2 \Rightarrow \min.$$



# Linear Regression



Model predicts output  $\mathbf{h}$  given input  $\mathbf{x}$

$$\mathbf{x} = \begin{bmatrix} x^{(1)} \\ \dots \\ x^{(m)} \end{bmatrix} \rightarrow \begin{bmatrix} 1 & x^{(1)} \\ \dots & \dots \\ 1 & x^{(m)} \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \dots \\ y^{(m)} \end{bmatrix}, \boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix}.$$

$$\mathbf{h} = \mathbf{x} \boldsymbol{\phi} = \begin{bmatrix} \phi_0 + \phi_1 x^{(1)} \\ \vdots \\ \phi_0 + \phi_1 x^{(m)} \end{bmatrix}$$

$$L(\boldsymbol{\phi}) = \frac{1}{2m} \sum_{i=1}^m (h^{(i)} - y^{(i)})^2 \Rightarrow \min.$$

$$\nabla L = \left[ \frac{\partial L}{\partial \phi_j} \right] = \frac{1}{m} \mathbf{x}^T (\mathbf{h} - \mathbf{y}).$$

Training algorithm.

1. Initialize the *weights*  $\boldsymbol{\phi}$  with a random seed

2. Calculate the *hypothesis* matrix  $\mathbf{h} = \mathbf{x}\boldsymbol{\phi}$  and the *loss gradient*:  $\nabla L = \frac{1}{m} \mathbf{x}^T (\mathbf{h} - \mathbf{y})$

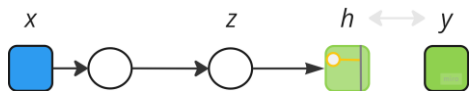
3. For the given  $\phi_j$  components (annotated with index 'prev',  $\phi_j^{prev}$ ) calculate the newer ones  $\phi_j^{next}$  moving towards the direction, which is opposite to the loss gradient vector, with steps which are proportional to the *learning rate*  $\alpha$ :

$$\boldsymbol{\phi}^{next} = \boldsymbol{\phi}^{prev} - \alpha \nabla L, \text{ or in the scalar form } \phi_0^{next} = \phi_0^{prev} - \alpha \frac{\partial L}{\partial \phi_0}, \phi_1^{next} = \phi_1^{prev} - \alpha \frac{\partial L}{\partial \phi_1}$$

4. Repeat pp. 2-3 until the minimum of the loss function  $L$  is reached, based on the condition of small changes in its value over several neighboring iterations or based on the condition of reaching the maximum number of iterations :  $L^{next} - L^{prev} < \delta$ , #iter. > max # of iter

5. Save the trained model (model weights):  $\boldsymbol{\phi}$ .

# Linear Regression. Generalization (multiple var., polynomial)



Model predicts output  $h$  given input  $x$

Consider a model  $f = [x^{(i)}, \phi]$  parameterized with weights  $\phi$  that maps each  $i$ -th input sample  $x^{(i)}$  into the output  $z^{(i)}$  which then transforms into the hypothesis  $h^{(i)}$  that should be close to the label  $y^{(i)}$ .

$$h = \phi_0 + \phi_1 x_1 + \dots + \phi_n x_n$$

$$x = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & \dots & x_n^{(m)} \end{bmatrix}; \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}; \quad \phi = \begin{bmatrix} \phi_0 \\ \phi_1 \\ \vdots \\ \phi_n \end{bmatrix}; \quad h = x\phi$$

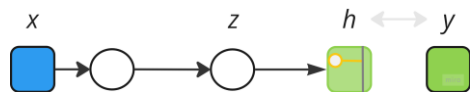
$$L(\phi) = \frac{1}{2m} \sum_{i=1}^m (h^{(i)} - y^{(i)})^2 \Rightarrow \min.$$

Training algorithm.

1. Initialize  $\phi$
2. Calculate  $h = x\phi$  and  $\nabla L = \frac{1}{m} x^T (h - y)$
3. Update  $\phi^{next} = \phi^{prev} - \alpha \nabla L$

4. Repeat pp. 2-3  $L^{next} - L^{prev} < \delta$ , #iter. > max # of iter
5. Save the trained model (model weights):  $\phi$ .

# Linear Regression. Generalization (multiple var., polynomial)



Model predicts output  $h$  given input  $x$

$$h = \phi_0 + \phi_1 x + \phi_2 x^2 + \dots + \phi_n x^n$$

$\downarrow$        $\downarrow$        $\downarrow$   
 $x_1$      $x_2$      $x_n$

$$h = \phi_0 + \phi_1 x_1 + \dots + \phi_n x_n \rightarrow \text{check the previous task}$$

Consider a model  $f = [x^{(i)}, \phi]$  parameterized with weights  $\phi$  that maps each  $i$ -th input sample  $x^{(i)}$  into the output  $z^{(i)}$  which then transforms into the hypothesis  $h^{(i)}$  that should be close to the label  $y^{(i)}$ .

$$L(\phi) = \frac{1}{2m} \sum_{i=1}^m (h^{(i)} - y^{(i)})^2 \Rightarrow \min.$$

Training algorithm.

1. Initialize  $\phi$
2. Calculate  $h = x\phi$  and  $\nabla L = \frac{1}{m} x^T (h - y)$
3. Update  $\phi^{next} = \phi^{prev} - \alpha \nabla L$
4. Repeat pp. 2-3  $L^{next} - L^{prev} < \delta$ , #iter. > max # of iter
5. Save the trained model (model weights):  $\phi$ .



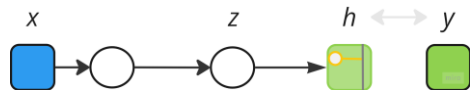
# Agenda

- I. Linear Regression and its Generalization
- II. Logistic Regression and its Generalization
- III. Setting of the models

All models are wrong, but some are useful.  
/George Box/



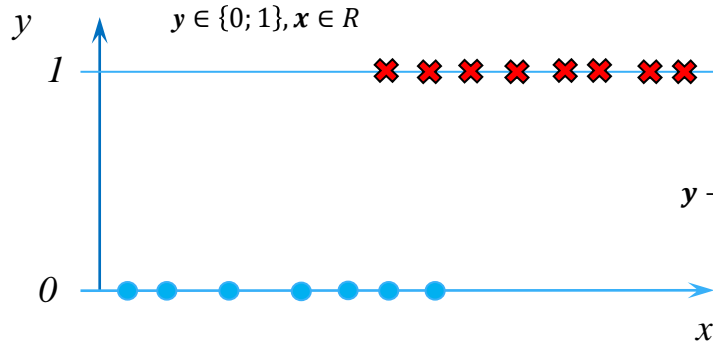
# Logistic Regression



Model predicts output  $h$  given input  $x$

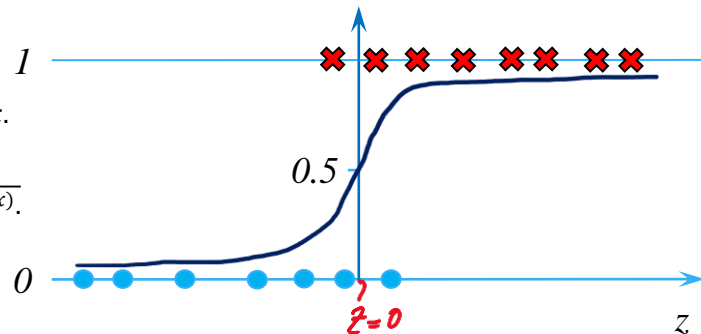
Consider a model  $f = [x^{(i)}, \phi]$  parameterized with weights  $\phi$  that maps each  $i$ -th input sample  $x^{(i)}$  into the output  $z^{(i)}$  which then transforms into the hypothesis  $h^{(i)}$  that should be close to the label  $y^{(i)}$ .

$$L(\phi) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \ln(h^{(i)}) + (1 - y^{(i)}) \ln(1 - h^{(i)})) \Rightarrow \min.$$



*transform*  
 $x \rightarrow z: z = \phi_0 + \phi_1 x.$

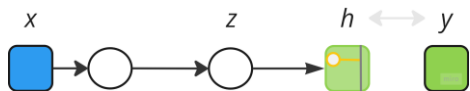
$$y \rightarrow h: h(z) = \frac{1}{1 + e^{-z(x)}}.$$



$$\mathbf{x} = \begin{bmatrix} x^{(1)} \\ \dots \\ x^{(m)} \end{bmatrix}; \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \dots \\ y^{(m)} \end{bmatrix}; \rightarrow \mathbf{x} = \begin{bmatrix} 1 & x^{(1)} \\ \dots & \dots \\ 1 & x^{(m)} \end{bmatrix}; \boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix}; \mathbf{z} = \mathbf{x}\boldsymbol{\phi};$$

$$\rightarrow h(z) = \frac{1}{1 + e^{-z(x)}}, \text{ or } \mathbf{h} = \sigma(\mathbf{z}).$$

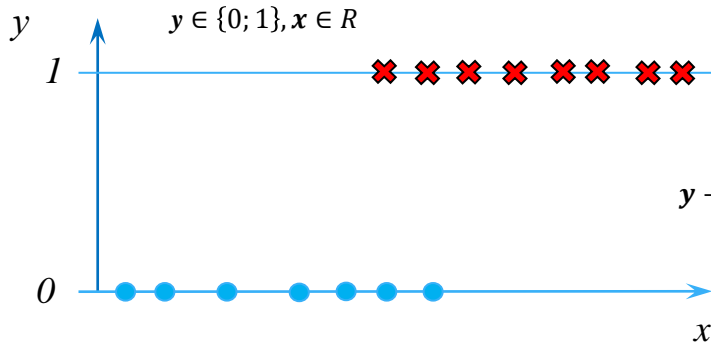
# Logistic Regression



Model predicts output  $h$  given input  $x$

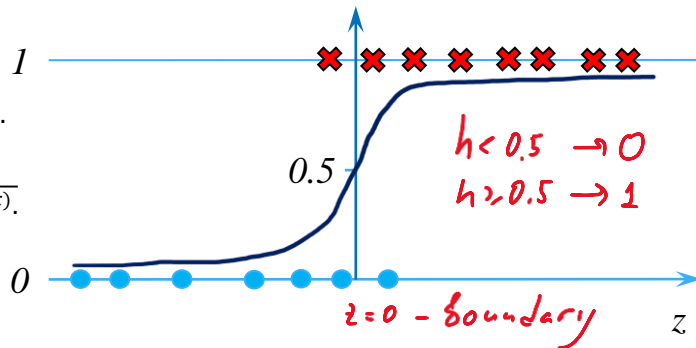
Consider a model  $f = [x^{(i)}, \phi]$  parameterized with weights  $\phi$  that maps each  $i$ -th input sample  $x^{(i)}$  into the output  $z^{(i)}$  which then transforms into the hypothesis  $h^{(i)}$  that should be close to the label  $y^{(i)}$ .

$$L(\phi) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \ln(h^{(i)}) + (1 - y^{(i)}) \ln(1 - h^{(i)})) \Rightarrow \min.$$



$$x \rightarrow z: z = \phi_0 + \phi_1 x.$$

$$y \rightarrow h: h(z) = \frac{1}{1 + e^{-z(x)}}.$$



Training algorithm.

1. Initialize  $\phi$

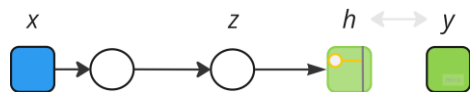
2. Calculate  $z = x\phi$ ,  $h = \sigma(z)$ , then  $\nabla L = \frac{1}{m} x^T (h - y)$

3. Update  $\phi^{next} = \phi^{prev} - \alpha \nabla L$

4. Repeat pp. 2-3  $L^{next} - L^{prev} < \delta$ , #iter. > max # of iter

5. Save the trained model (model weights):  $\phi$ .

# Logistic Regression. Generalization (multiple var., polynomial)



Model predicts output  $h$  given input  $x$

$$z = \phi_0 + \phi_1 x_1 + \phi_2 x_2, \quad h = \sigma(z).$$

Consider a model  $f = [x^{(i)}, \phi]$  parameterized with weights  $\phi$  that maps each  $i$ -th input sample  $x^{(i)}$  into the output  $z^{(i)}$  which then transforms into the hypothesis  $h^{(i)}$  that should be close to the label  $y^{(i)}$ .

$$L(\phi) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \ln(h^{(i)}) + (1 - y^{(i)}) \ln(1 - h^{(i)})) \Rightarrow \min.$$

$$x = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} \\ \dots & \dots \\ x_1^{(m)} & x_2^{(m)} \end{bmatrix}; y = \begin{bmatrix} y^{(1)} \\ \dots \\ y^{(m)} \end{bmatrix}; \quad \rightarrow \quad x = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ \dots & \dots & \dots \\ 1 & x_1^{(m)} & x_2^{(m)} \end{bmatrix}; \phi = \begin{bmatrix} \phi_0 \\ \phi_1 \\ \phi_2 \end{bmatrix}; z = x\phi; \quad \rightarrow \quad h(z) = \frac{1}{1 + e^{-z(x)}}, \text{ or } h = \sigma(z).$$

$$\text{boundary: } 0 = \phi_0 + \phi_1 x_1 + \phi_2 x_2$$

$$x_2 = -\frac{\phi_0}{\phi_2} - \frac{\phi_1}{\phi_2} x_1$$

Training algorithm.

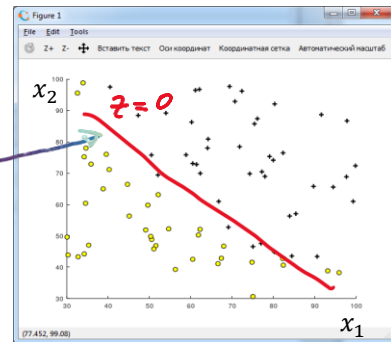
1. Initialize  $\phi$

2. Calculate  $z = x\phi$ ,  $h = \sigma(z)$ , then  $\nabla L = \frac{1}{m} x^T (h - y)$

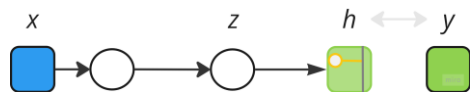
3. Update  $\phi^{next} = \phi^{prev} - \alpha \nabla L$

4. Repeat pp. 2-3  $L^{next} - L^{prev} < \delta$ , #iter. > max # of iter

5. Save the trained model (model weights):  $\phi$ .



# Logistic Regression. Generalization (multiple var., polynomial)



Model predicts output  $h$  given input  $x$



[Gray scale picture of "Nine"](#)

Consider a model  $f = [x^{(i)}, \phi]$  parameterized with weights  $\phi$  that maps each  $i$ -th input sample  $x^{(i)}$  into the output  $z^{(i)}$  which then transforms into the hypothesis  $h^{(i)}$  that should be close to the label  $y^{(i)}$ .

$$L(\phi) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \ln(h^{(i)}) + (1 - y^{(i)}) \ln(1 - h^{(i)})) \Rightarrow \min.$$

$$h = \phi_0 + \phi_1 x_1 + \dots + \phi_n x_n$$

$$\mathbf{x} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix}; \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix};$$

$$\rightarrow \mathbf{x} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix}; \phi = \begin{bmatrix} \phi_0 \\ \vdots \\ \phi_n \end{bmatrix}; \mathbf{z} = \mathbf{x}\phi; \mathbf{h} = \sigma(\mathbf{z}).$$

Training algorithm.

1. Initialize  $\phi$

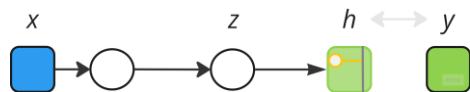
2. Calculate  $\mathbf{z} = \mathbf{x}\phi$ ,  $\mathbf{h} = \sigma(\mathbf{z})$ , then  $\nabla L = \frac{1}{m} \mathbf{x}^T (\mathbf{h} - \mathbf{y})$

3. Update  $\phi^{next} = \phi^{prev} - \alpha \nabla L$

4. Repeat pp. 2-3  $L^{next} - L^{prev} < \delta$ , #iter. > max # of iter

5. Save the trained model (model weights):  $\phi$ .

# Logistic Regression. Generalization (multiple var., polynomial)



Model predicts output  $h$  given input  $x$

$$z = \phi_0 + \phi_1 x_1 + \phi_2 x_2$$

$\downarrow$   $x_1$        $\downarrow$   $x_2$

$$\mathbf{x} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} \\ \dots & \dots \\ x_1^{(m)} & x_2^{(m)} \end{bmatrix}; \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \dots \\ y^{(m)} \end{bmatrix}; \rightarrow$$

$$z = \phi_0 + \phi_1 x_1 + \phi_2 x_2 - \text{Check the previous task.}$$

Consider a model  $f = [x^{(i)}, \phi]$  parameterized with weights  $\phi$  that maps each  $i$ -th input sample  $x^{(i)}$  into the output  $z^{(i)}$  which then transforms into the hypothesis  $h^{(i)}$  that should be close to the label  $y^{(i)}$ .

$$L(\phi) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \ln(h^{(i)}) + (1 - y^{(i)}) \ln(1 - h^{(i)})) \Rightarrow \min.$$

more general:

$$z = \phi_0 + \phi_1 x_1 + \dots + \phi_n x_1^n + \phi_{n+1} x_2 + \dots + \phi_k x_2^2$$

Training algorithm.

1. Initialize  $\phi$

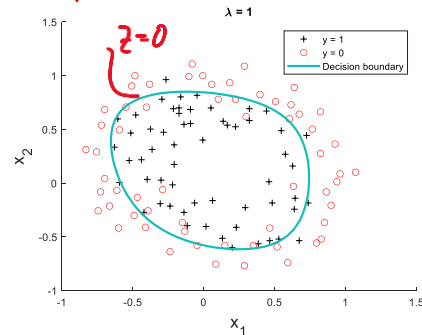
2. Calculate  $\mathbf{z} = \mathbf{x}\phi$ ,  $\mathbf{h} = \sigma(\mathbf{z})$ , then  $\nabla L = \frac{1}{m} \mathbf{x}^T (\mathbf{h} - \mathbf{y})$

3. Update  $\phi^{next} = \phi^{prev} - \alpha \nabla L$

4. Repeat pp. 2-3  $L^{next} - L^{prev} < \delta$ ,

#iter. > max # of iter

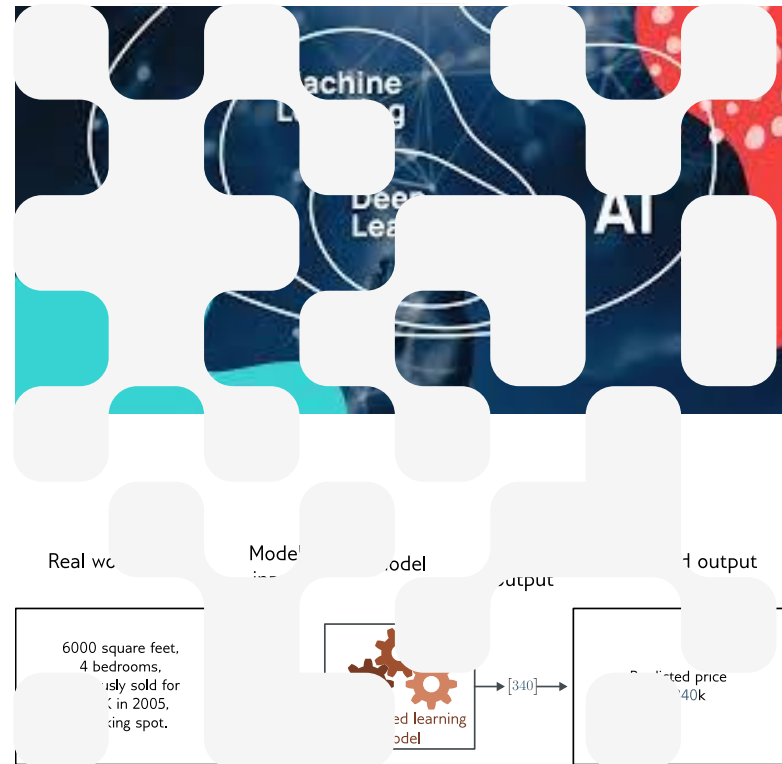
5. Save the trained model (model weights):  $\phi$ .



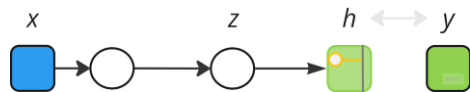
# Agenda

- I. Linear Regression and its Generalization
- II. Logistic Regression and its Generalization
- III. Setting of the models

All models are wrong, but some are useful.  
/George Box/



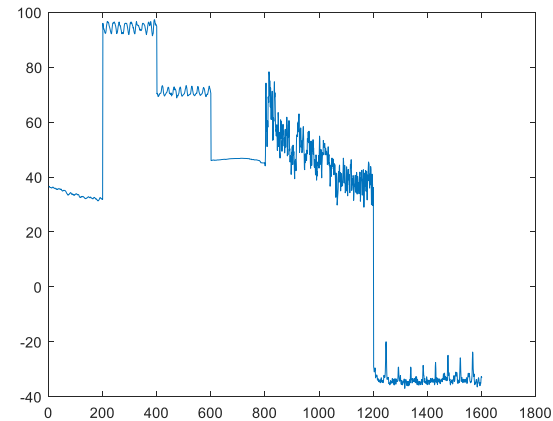
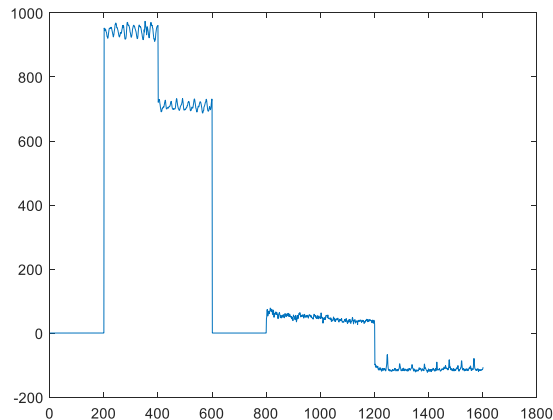
# ML Settings



Model predicts output  $h$  given input  $x$

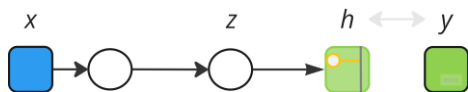
Model *parameters* are determined during the solution of the ML problem. For example, in regression problems, the parameters are the components of the matrix of weights  $\phi$ . *Hyperparameters* are set by the user, usually not in a single way, and their values affect the values of the sought parameters.

## 1. Feature Scaling





# ML Settings



Model predicts output  $h$  given input  $x$

1. Feature Scaling
2. Learning Rate
3. Error and # of iterations
4. **Regularization (L2)**

Model *parameters* are determined during the solution of the ML problem. For example, in regression problems, the parameters are the components of the matrix of weights  $\phi$ . *Hyperparameters* are set by the user, usually not in a single way, and their values affect the values of the sought parameters.

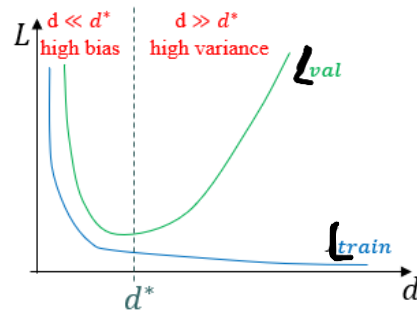
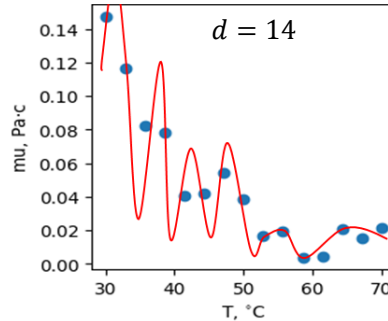
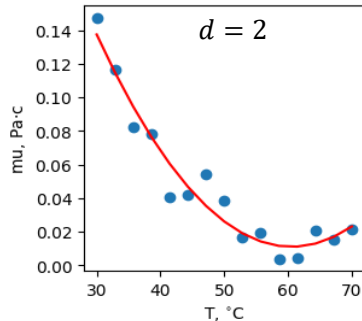
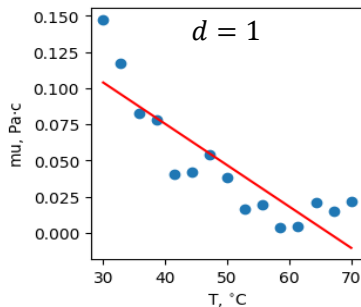
Training  $\{(x_i, y_i)\}$

validation

test

$$L = \frac{1}{2m} \left[ \sum_{i=1}^m (h^{(i)} - y^{(i)})^2 + \lambda \sum_{j=1}^n \phi_j^2 \right] \Rightarrow \min.$$

$$h(x) = \theta_j x^j, (j = 0, \dots, d)$$



## Just think about it



1. How can the gradient descent method be improved to find global minima instead of local ones?
2. Can the discussed linear regression problems be solved analytically without using the gradient descent method?
3. Why is the use of high-degree polynomials generally not recommended when building regression models?

Thank you for your attention!

[a.kornaev@innopolis.ru](mailto:a.kornaev@innopolis.ru), [@avkornaev](#)











