

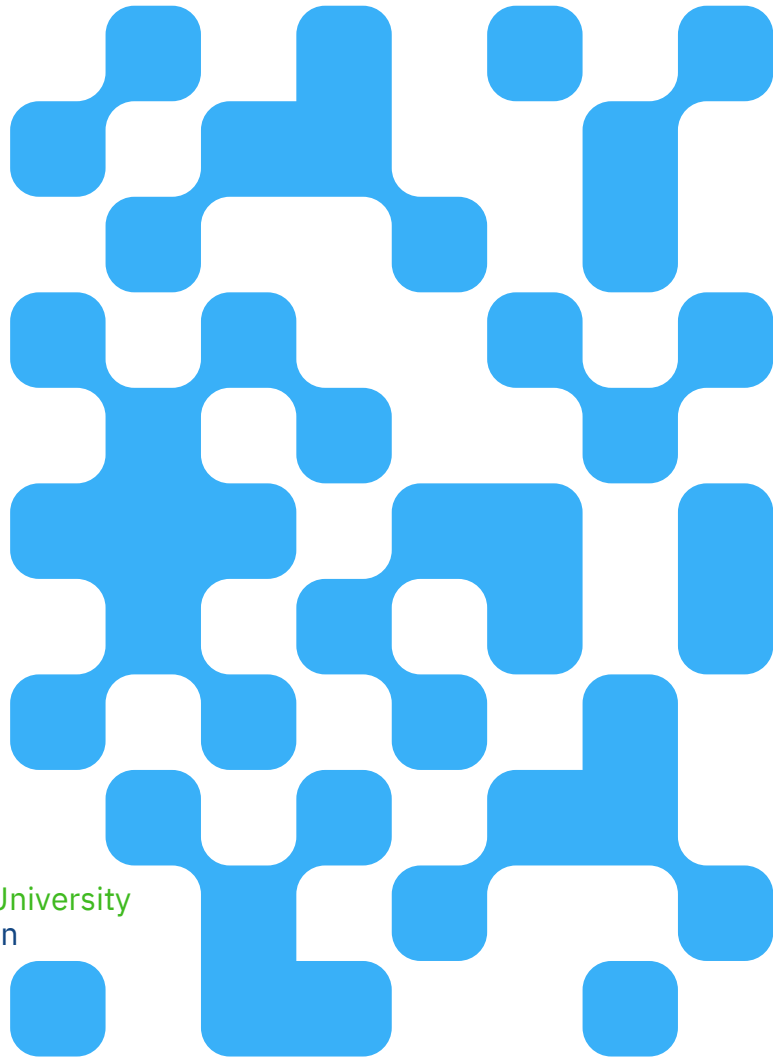


Machine Learning

2025 (ML-2025)

Lecture 2. Linear models

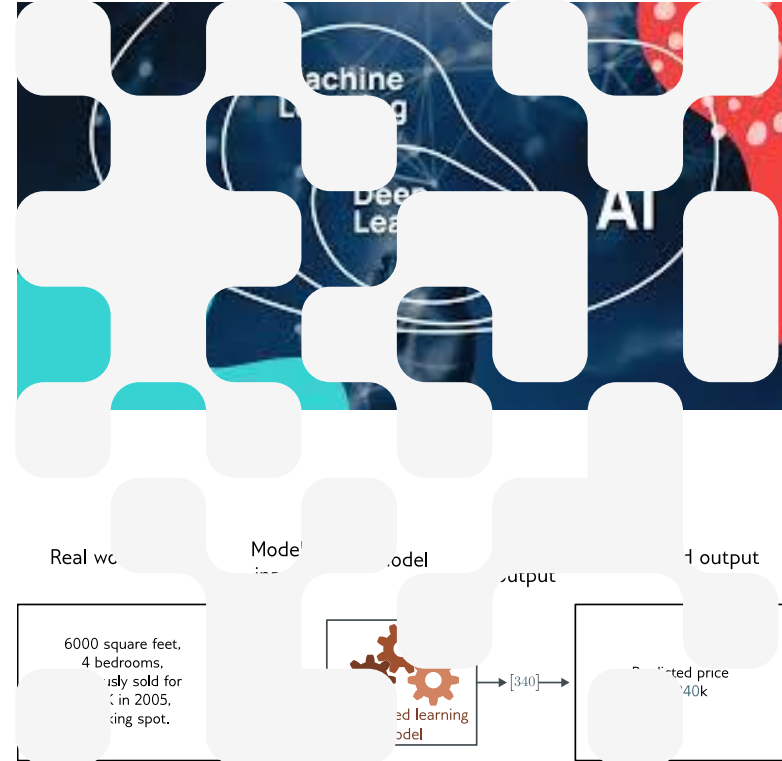
by Alexei Kornaev, Dr. Sc., Assoc. Prof., Robotics and CV, [Innopolis University](#)
Researcher at the RC for AI, [National RC for Oncology n.a. NN Blokhin](#)



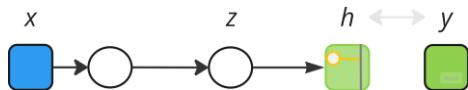
Agenda

- 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/



Linear Regression



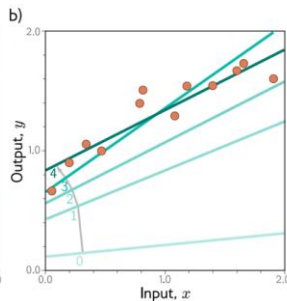
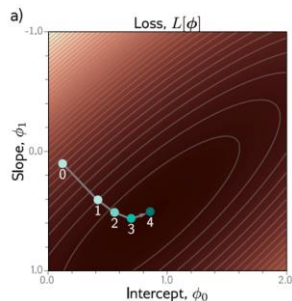
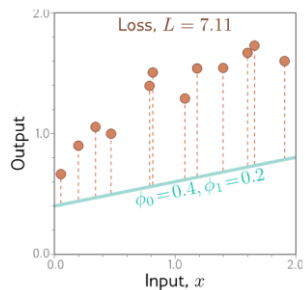
Model predicts output h given input x

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

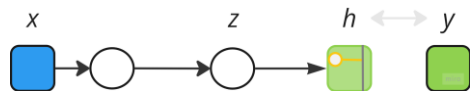
$$\mathbf{x} = \begin{bmatrix} x^{(1)} \\ \dots \\ x^{(m)} \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \dots \\ y^{(m)} \end{bmatrix}, \boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix}.$$

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}

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

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

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

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 ϕ_j components (annotated with index 'prev', ϕ_j^{prev}) calculate the newer ones ϕ_j^{next} moving towards the direction, which is opposite to the loss gradient vector, with steps which are proportional to the *learning rate* α :

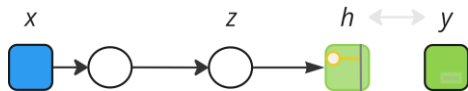
$$\boldsymbol{\phi}^{next} = \boldsymbol{\phi}^{prev} - \alpha \nabla L, \text{ or in the scalar form } \phi_0^{next} = \phi_0^{prev} - \alpha \frac{\partial L}{\partial \phi_0}, \phi_1^{next} = \phi_1^{prev} - \alpha \frac{\partial L}{\partial \phi_1}$$

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): $\boldsymbol{\phi}$.

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

Linear Regression. Generalization (multiple var., polynomial)



Model predicts output h given input x

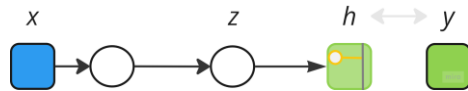
Consider a model $f = [x^{(i)}, \phi]$ parameterized with weights ϕ that maps each i -th input sample $x^{(i)}$ into the output $z^{(i)}$ which then transforms into the hypothesis $h^{(i)}$ that should be close to the label $y^{(i)}$.

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

Training algorithm.

1. Initialize ϕ
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): ϕ .

Linear Regression. Generalization (multiple var., polynomial)



Model predicts output h given input x

Consider a model $f = [x^{(i)}, \phi]$ parameterized with weights ϕ 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 ϕ
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): ϕ .

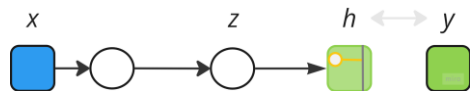
Agenda

- 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/



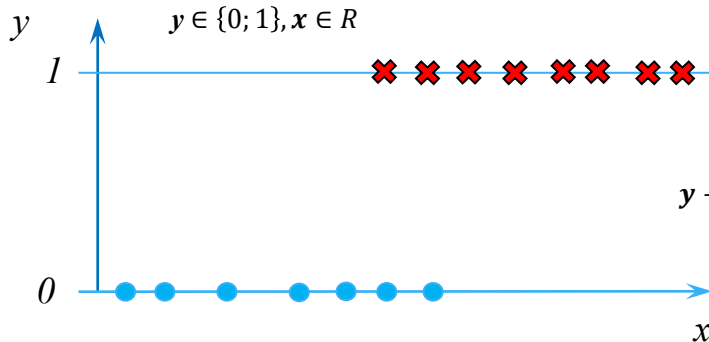
Logistic Regression



Model predicts output h given input x

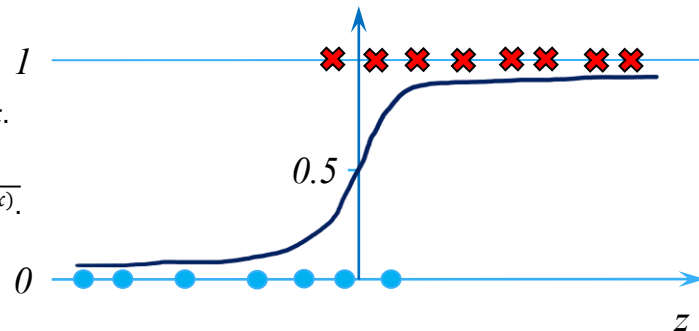
Consider a model $f = [x^{(i)}, \phi]$ parameterized with weights ϕ that maps each i -th input sample $x^{(i)}$ into the output $z^{(i)}$ which then transforms into the hypothesis $h^{(i)}$ that should be close to the label $y^{(i)}$.

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



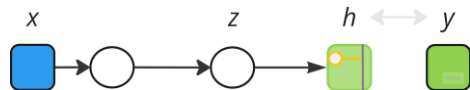
$$x \rightarrow z: z = \phi_0 + \phi_1 x.$$

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



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

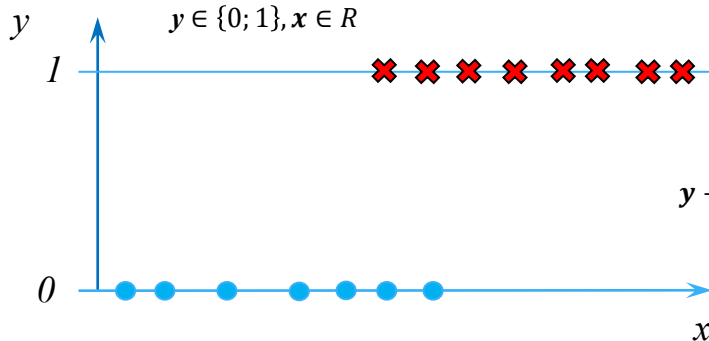
Logistic Regression



Model predicts output h given input x

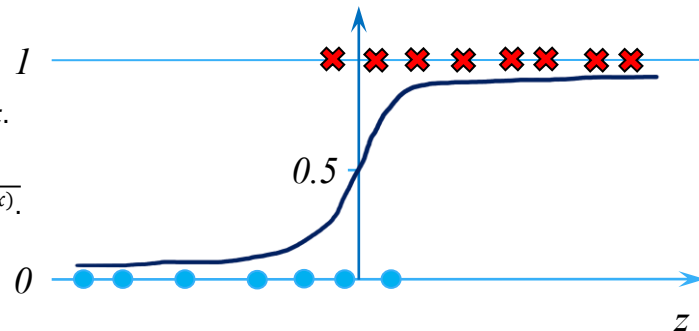
Consider a model $f = [x^{(i)}, \phi]$ parameterized with weights ϕ that maps each i -th input sample $x^{(i)}$ into the output $z^{(i)}$ which then transforms into the hypothesis $h^{(i)}$ that should be close to the label $y^{(i)}$.

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



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

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



Training algorithm.

1. Initialize ϕ

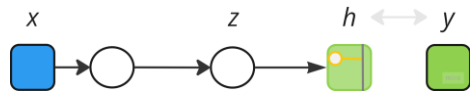
2. Calculate $z = x\phi$, $h = \sigma(z)$, then $\nabla L = \frac{1}{m} x^T (h - y)$

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): ϕ .

Logistic Regression. Generalization (multiple var., polynomial)



Model predicts output h given input x

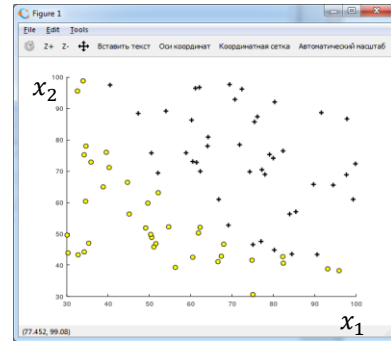
Consider a model $f = [x^{(i)}, \phi]$ parameterized with weights ϕ that maps each i -th input sample $x^{(i)}$ into the output $z^{(i)}$ which then transforms into the hypothesis $h^{(i)}$ that should be close to the label $y^{(i)}$.

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

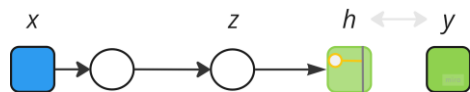
$$\mathbf{x} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} \\ \vdots & \vdots \\ x_1^{(m)} & x_2^{(m)} \end{bmatrix}; \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}; \quad \rightarrow \quad \mathbf{x} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ \vdots & \vdots & \vdots \\ 1 & x_1^{(m)}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

Training algorithm.

1. Initialize ϕ
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): ϕ .



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 ϕ 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 $y^{(i)}$.

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

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

Training algorithm.

1. Initialize ϕ

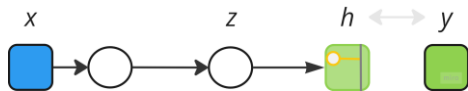
2. Calculate $\mathbf{z} = \mathbf{x}\phi$, $\mathbf{h} = \sigma(\mathbf{z})$, then $\nabla L = \frac{1}{m} \mathbf{x}^T (\mathbf{h} - \mathbf{y})$

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): ϕ .

Logistic Regression. Generalization (multiple var., polynomial)



Model predicts output h given input x

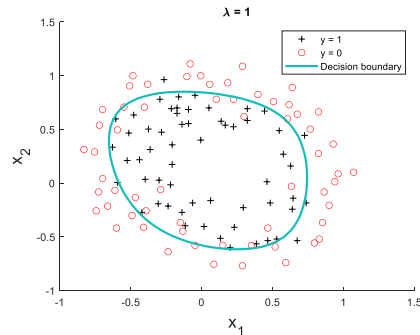
Consider a model $f = [x^{(i)}, \phi]$ parameterized with weights ϕ that maps each i -th input sample $x^{(i)}$ into the output $z^{(i)}$ which then transforms into the hypothesis $h^{(i)}$ that should be close to the label $y^{(i)}$.

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

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

Training algorithm.

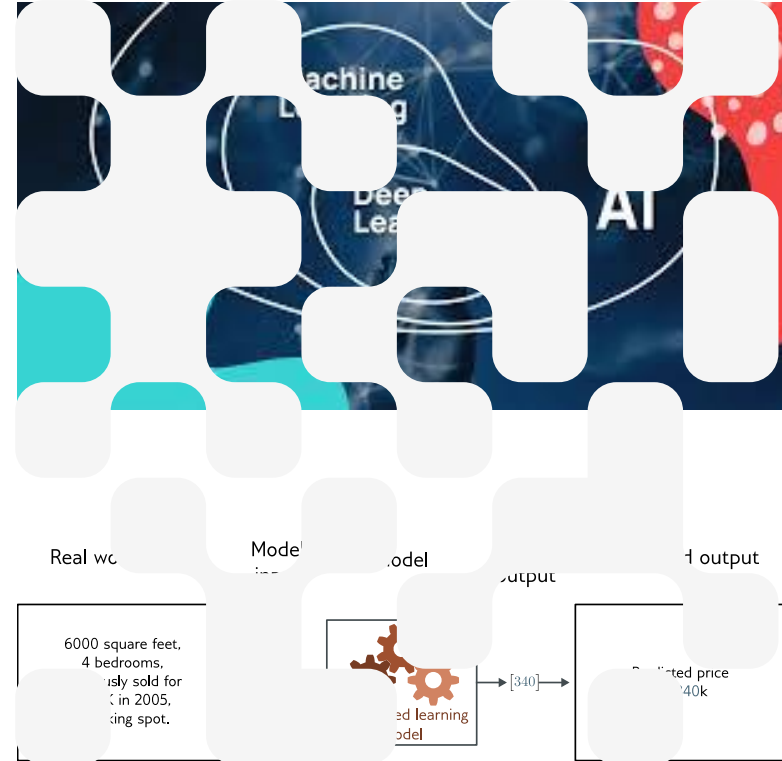
1. Initialize ϕ
2. Calculate $\mathbf{z} = \mathbf{x}\phi$, $\mathbf{h} = \sigma(\mathbf{z})$, then $\nabla L = \frac{1}{m} \mathbf{x}^T (\mathbf{h} - \mathbf{y})$
3. Update $\phi^{next} = \phi^{prev} - \alpha \nabla L$
4. Repeat pp. 2-3 $L^{next} - L^{prev} < \delta$,
#iter. > max # of iter
5. Save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



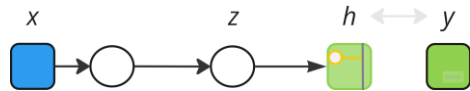
Agenda

- 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/



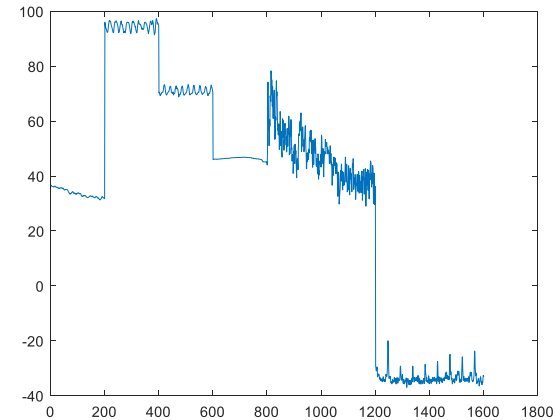
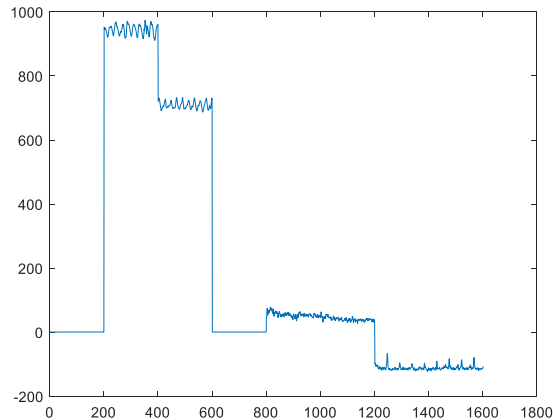
ML Settings



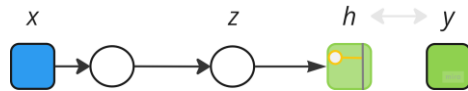
Model predicts output h given input x

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 ϕ . *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

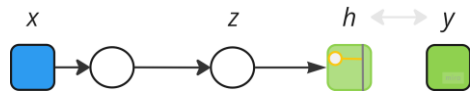


Model predicts output h given input x

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 ϕ . *Hyperparameters* are set by the user, usually not in a single way, and their values affect the values of the sought parameters.

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 ϕ . *Hyperparameters* are set by the user, usually not in a single way, and their values affect the values of the sought parameters.

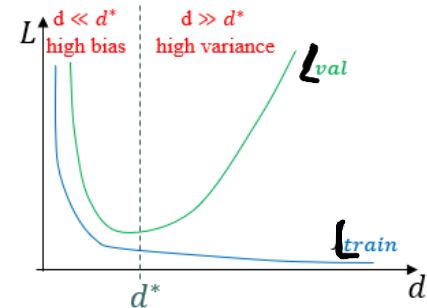
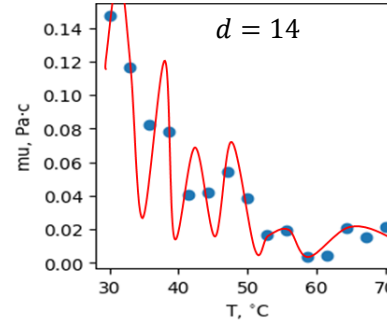
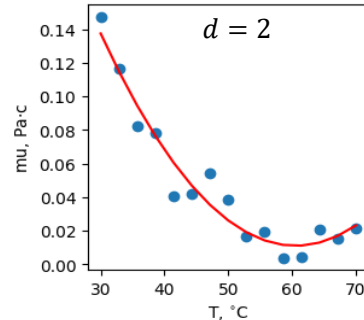
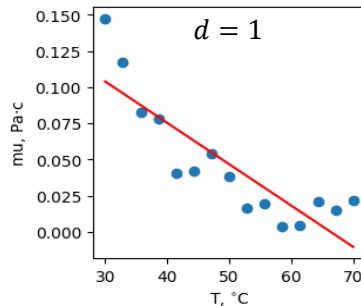
Training $\{(x_i, y_i)\}$

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, [@avkornaev](#)



