



# Feature Ranking via Eigenvector Centrality

Team 35:

- Ivan Baybuza
- Aleksandr Belov
- Mariia Kopylova
- Miron Kuznetsov

# Introduction

- Dataset: 165k clinical records by 55k patients
- Source of data: Moscow Cardiology center
- Records parsed to 88 features (used OHE, text processing)
- $1.6 \times 10^7$  values in matrix X
- Matrix is sparse, only 13% are not NaN

## Problem statement

- Target variable is “fibrillation” (“фибрилляция предсердий”)
- Trying to find predictors of this diagnosis

→ Project by Dmitrii Dylov with Moscow Cardiology Center

→ Doctors have many different scales (i.e. on the right)

Число присваиваемых баллов	Признак
(максимальный суммарный балл для каждого пациента – 9)	
1 балл	Женский пол
1 балл	Возраст от 65 до 74 лет
2 балла	Возраст 75 лет и старше
1 балл	Артериальная гипертония

# Main idea from paper <sup>[1]</sup>

Graph-based method for feature selection:

$$X = \{x^{(1)}, \dots, x^{(n)}\}$$

$$G = (V, E) \quad \longrightarrow \quad \text{Adjacency matrix } A$$

$$a_{ij} = \varphi(x^{(i)}, x^{(j)})$$

# ECFS algorithm

$$f_i = \frac{|\mu_{i,1} - \mu_{i,2}|^2}{\sigma_{i,1}^2 + \sigma_{i,2}^2}$$

$$m_i = \sum_{y \in Y} \sum_{z \in x^{(i)}} p(z, y) \log \left( \frac{p(z, y)}{p(z)p(y)} \right)$$

$$\Sigma(i, j) = \max \left( \sigma^{(i)}, \sigma^{(j)} \right)$$

$$k = (f \cdot m^\top)$$

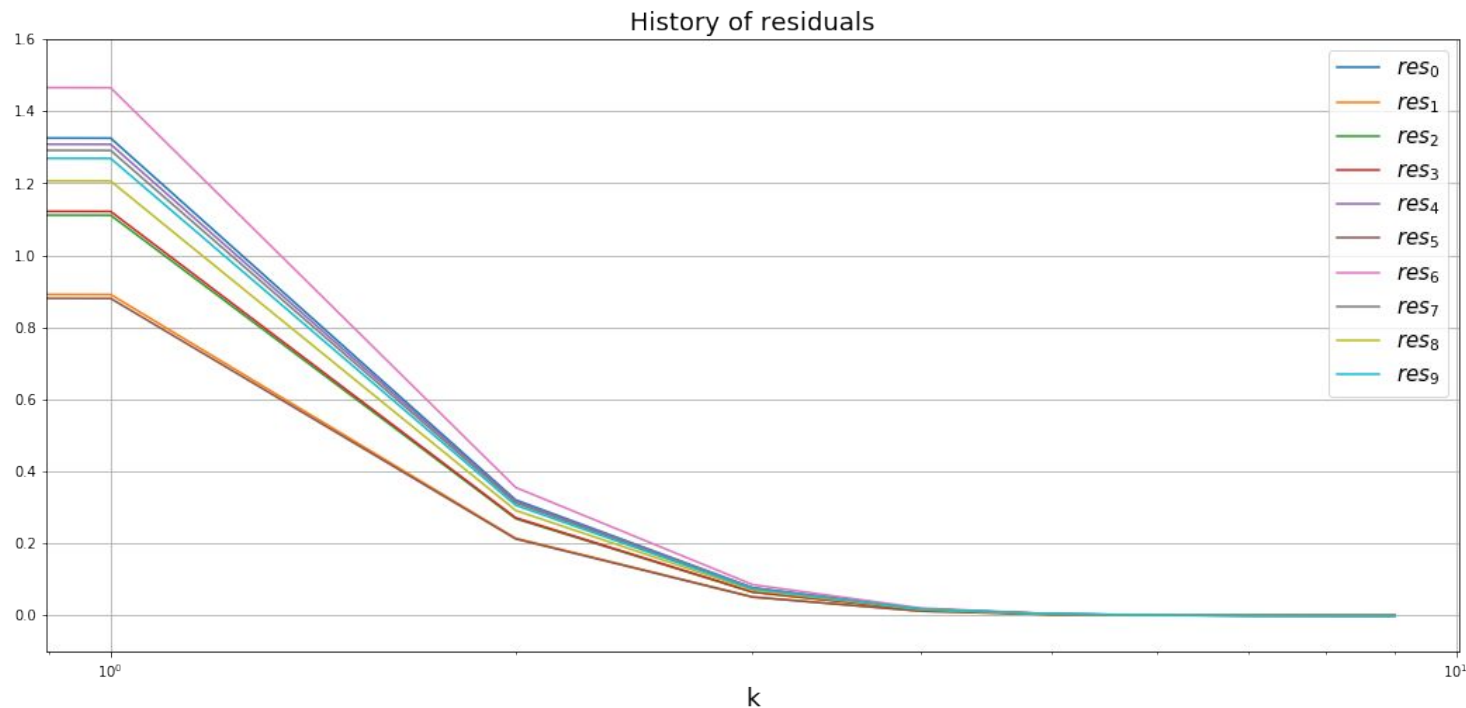
$$A = \alpha k + (1 - \alpha) \Sigma$$

where  $a_{ij} = \varphi(x^{(i)}, x^{(j)})$  is a pairwise potential term

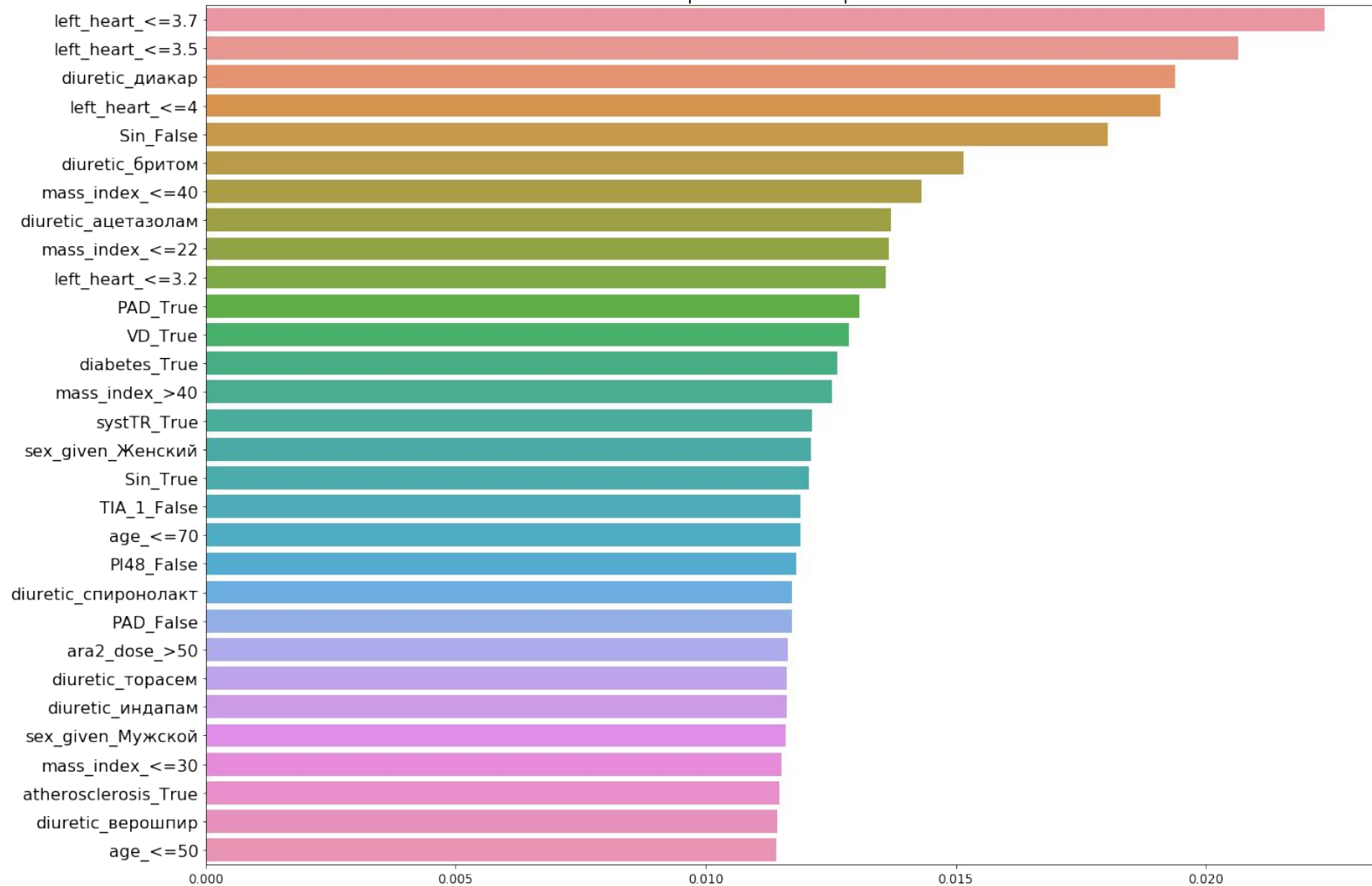
 Use the Power method !!!

$$A^n e \rightarrow v_1$$

10 iterations is  
more than enough



Top-30 feature impotance



# Obtained results

The main predictors of Atrial fibrillation is:

- 1) Large size of the left heart ( $> 3.2$ )
- 2) Large Mass index of a patient (*2 degree of obesity and more*)
- 3) An absence of Sinus tachycardia
- 4) A presence of diabetes



# Conclusion

- The investigated *Eigenvector Centrality* method may be applicable to this type of problems and could give some unobvious information
- The results should be delivered to the specialists in Cardiology Center to be considered by them

# References

1. Ranking to Learn: Feature Ranking and Selection via Eigenvector Centrality by Giorgio Roffo and Simone Melzi (2017)
2. Mining electronic medical records to explore the linkage between healthcare resource utilization and disease severity in diabetic patients
3. Entity Embeddings of Categorical Variables by Cheng Guo and Felix Berkhahn (2016)
4. [https://github.com/albellov/nla\\_project](https://github.com/albellov/nla_project)