# Assignment for *Bioinformatics aspects of aging and rejuvenation* lecture

Zoltán Szarvas[1]

lecturer: Csaba Kerepesi[2]

[1]Faculty of Computer Science, Eötvös Loránd University
[2]HUN-REN Institute for Computer Science and Control

May 2024

# Assignment 1

Train epigenetic clock for
microarray-based methylation dataset[3]

# Dataset

## Source

- Genome-wide Methylation Profiles Reveal Quantitative Views of Human Aging Rates (Gene Expression Omnibus) [5] [6]
- approximately 450k CpGs from **human whole blood**
- Illumina Infinium 450k (microarray)
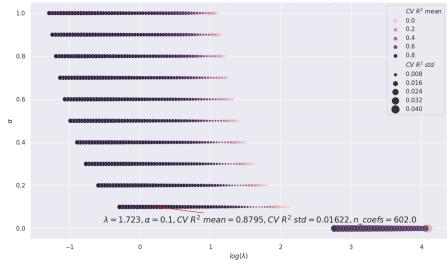
## Details

- sample N=656
- features=473 034

## Model

- Elasticnet using Glmnet [1]
- 80-20 train-test split
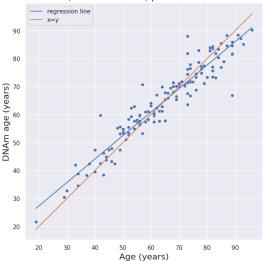- 10-fold CV

# Hyperparameter Optimization (ElasticNet)



Hyperparameter optimization result on training data
best is selected based on $(CV\ R^2\ mean) - 0.5 \cdot (CV\ R^2\ std)$

$\lambda = 1.723, \alpha = 0.1, CV\ R^2\ mean = 0.8795, CV\ R^2\ std = 0.01622, n\_coefs = 602.0$

# Results on Test Set



Test set ($n = 132$); $\alpha = 0.1$, $\lambda = 1.72$
$R^2 = 0.903$, stderr=0.0242, p=1.1e-67 $MedAE = 2.83$

# Assignment 2

Train epigenetic clock for RRBS-based methylation dataset (with missing data)[4]

# Dataset

## Source

- Data from "Depression and suicide risk prediction models using blood-derived multi-omics data" [2]
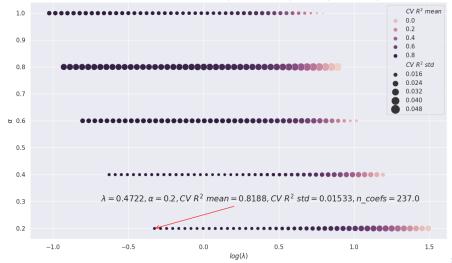- Reduced Representation Bisulfite Sequencing (RRBS)

## Details

- sample N=182 (Healthy/Control=87, Major Depressive Disorder (MDD)=39, Suicide Attempters (SA)=56)
- features=8 722 096

## Model

- Elasticnet using Glmnet [1]
- 80-20 train-test split, 10-fold CV
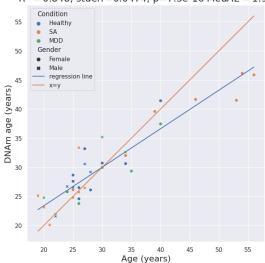- data for missing CpG sites was filled with **mean**

# Hyperparameter Optimization (ElasticNet)



Hyperparameter optimization result on training data
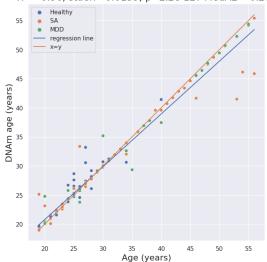best is selected based on $(CV\ R^2\ mean) - 0.5 \cdot (CV\ R^2\ std)$

$\lambda = 0.4722, \alpha = 0.2, CV\ R^2\ mean = 0.8188, CV\ R^2\ std = 0.01533, n\_coefs = 237.0$

# Results on Test Set



Test set ($n = 37$); $\alpha = 0.2$, $\lambda = 0.472$
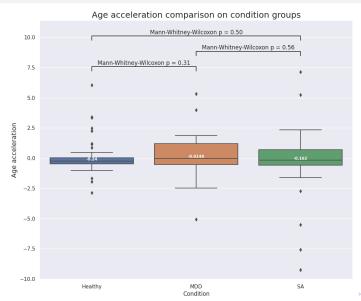$R^2 = 0.848$, stderr=0.0474, p=7.3e-16 $MedAE = 1.98$

# Results for All Instances



All samples ($n = 182$); $\alpha = 0.2$, $\lambda = 0.472$
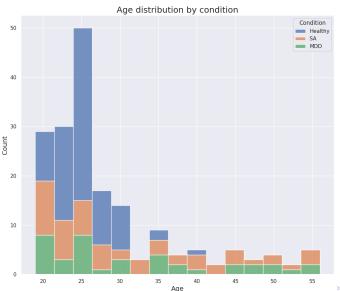$R^2 = 0.96$, stderr=0.0139, p=2.2e-127 $MedAE = 0.201$
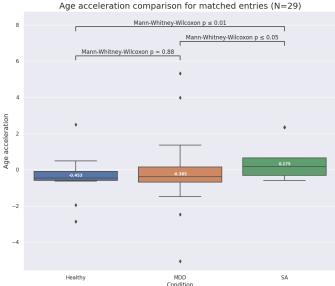
# Age acceleration comparison by condition



Age acceleration comparison on condition groups

# Age distribution by condition



Age distribution by condition

# Age acceleration comparison by condition (controlled)



Age acceleration comparison for matched entries (N=29)

# References I

[1] *An Introduction to 'glmnet' — glmnet.stanford.edu*.
https://glmnet.stanford.edu/articles/glmnet.html.
[Accessed 16-05-2024].

[2] Youngjune Bhak et al. "Depression and suicide risk prediction models using blood-derived multi-omics data". en. In: *Transl. Psychiatry* 9.1 (Oct. 2019), p. 262.

[3] *Epigenetic clock assignment 1*.
https://github.com/szazo/epigenetic-clock/blob/main/assignment1_microarray.ipynb.
2024.

[4] *Epigenetic clock assignment 2*.
https://github.com/szazo/epigenetic-clock/blob/main/assignment2_rrbs.ipynb. 2024.

## References II

[5] *Genome-wide Methylation Profiles Reveal Quantitative Views of Human Aging Rates*. https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE40279. [Accessed 13-05-2024]. 2012.

[6] Gregory Hannum et al. "Genome-wide Methylation Profiles Reveal Quantitative Views of Human Aging Rates". In: *Molecular Cell* 49.2 (Jan. 2013), pp. 359–367. ISSN: 1097-2765. DOI: 10.1016/j.molcel.2012.10.016. URL: http://dx.doi.org/10.1016/j.molcel.2012.10.016.