



HAL
open science

Impact of joint Dimension Reduction methods for survival prediction - Extension of a multi-omics benchmark study

Vincent Le Goff, Vincent Guillemot, Cathy Philippe, Gwendoline Mendes, Jean-François Deleuze, Edith Le Floch, Arnaud Gloaguen

► To cite this version:

Vincent Le Goff, Vincent Guillemot, Cathy Philippe, Gwendoline Mendes, Jean-François Deleuze, et al.. Impact of joint Dimension Reduction methods for survival prediction - Extension of a multi-omics benchmark study. Journées Ouvertes en Biologie, Informatique et Mathématiques (JOBIM), Jun 2024, Toulouse, France. cea-04734003

HAL Id: cea-04734003

<https://cea.hal.science/cea-04734003v1>

Submitted on 13 Oct 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Vincent LE GOFF¹, Vincent GUILLEMOT², Cathy PHILIPPE³, Gwendoline MENDES⁴, Jean-François DELEUZE¹, Edith LE FLOCH¹, Arnaud GLOAGUEN¹
¹ Centre National de Recherche en Génomique Humaine (CNRGH), CEA, Université Paris-Saclay, Evry, France

² Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, Paris, France

³ NeuroSpin, CEA Saclay, Université Paris-Saclay, France

⁴ CentraleSupélec, Gif-sur-Yvette, France

@vincent.legoff@cnrgh.fr

Introduction

This work aims to expand on a previous study [1] that compares the **survival predictions** of **13 different supervised multi-omics methods** on **datasets from 18 cancer types** from The Cancer Genome Atlas (TCGA).

This comparison is based on the analysis of **4 types of molecular data** (mRNA, miRNA, CNV, Mutations) and **clinical data**. Methods are divided into 3 groups: **reference models** (Kaplan-Meier and a Cox regression on clinical data only), **"naive" models**, unable to distinguish between each omics data, and **"structured" models**.

Objectives:

- Building on the conclusion in [1] that structured methods perform better than naive ones, we include methods not only capable of exploiting the group structure of omics data but also of extracting links between them: **Joint Dimension Reduction (JDR)** methods. We test these methods in **unsupervised and supervised** settings whenever possible.
- Then, to go further on the evaluation of the importance of clinical data **we compare all of these methods (naive/structured/JDR) without clinical data** to see if molecular data can lead to similar prediction performance as with the addition of clinical data.

JDR methods

RGCCA/SGCCA

Regularized Generalized Canonical Correlation Analysis [2]

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{w}_1, \dots, \mathbf{w}_L} \sum_{k,l=1}^L c_{kl} g(\operatorname{cov}(\mathbf{X}_k \mathbf{w}_k, \mathbf{X}_l \mathbf{w}_l)) \\ & \text{s.t.} \begin{cases} (1-\tau_l) \operatorname{var}(\mathbf{X}_l \mathbf{w}_l) + \tau_l \|\mathbf{w}_l\|_2^2 = 1, l = 1, \dots, L \\ \|\mathbf{w}_l\|_2^2 = 1 \text{ and } \|\mathbf{w}_l\|_1 \leq s_l \end{cases} \end{aligned}$$

JIVE

Joint and Individual Variation Explained [3]

$$\begin{aligned} & \mathbf{X}_l^T \approx \mathbf{U}_l \mathbf{S}_l + \mathbf{W}_l \mathbf{S}_l \text{ s.t.} \begin{cases} \mathbf{S}_l \mathbf{S}_l^T = 0 \\ \operatorname{rank}(\mathbf{S}_l) = r_j \\ \operatorname{rank}(\mathbf{S}_l) = r_a \end{cases} \end{aligned}$$

IntNMF

Integrative Non-Negative Matrix Factorization [4]

$$\mathbf{X}_l \approx \mathbf{W}_l \mathbf{H}_l^T \text{ s.t. } \mathbf{W}_l, \mathbf{H}_l \geq 0, l = 1, \dots, L$$

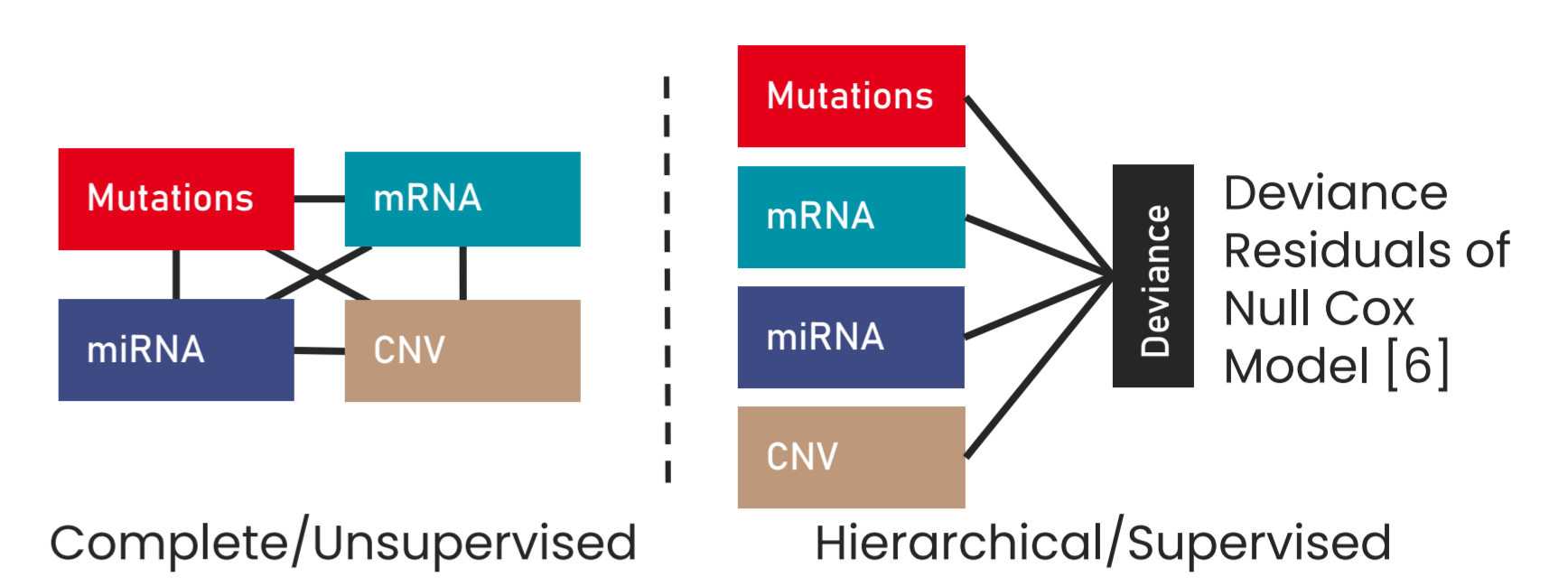
MOFA

Multi-Omics Factor Analysis [5]

$$\mathbf{X}_l \approx \mathbf{Z} \mathbf{W}_l^T$$

fit through Bayesian methods

RGCCA Unsupervised & Supervised



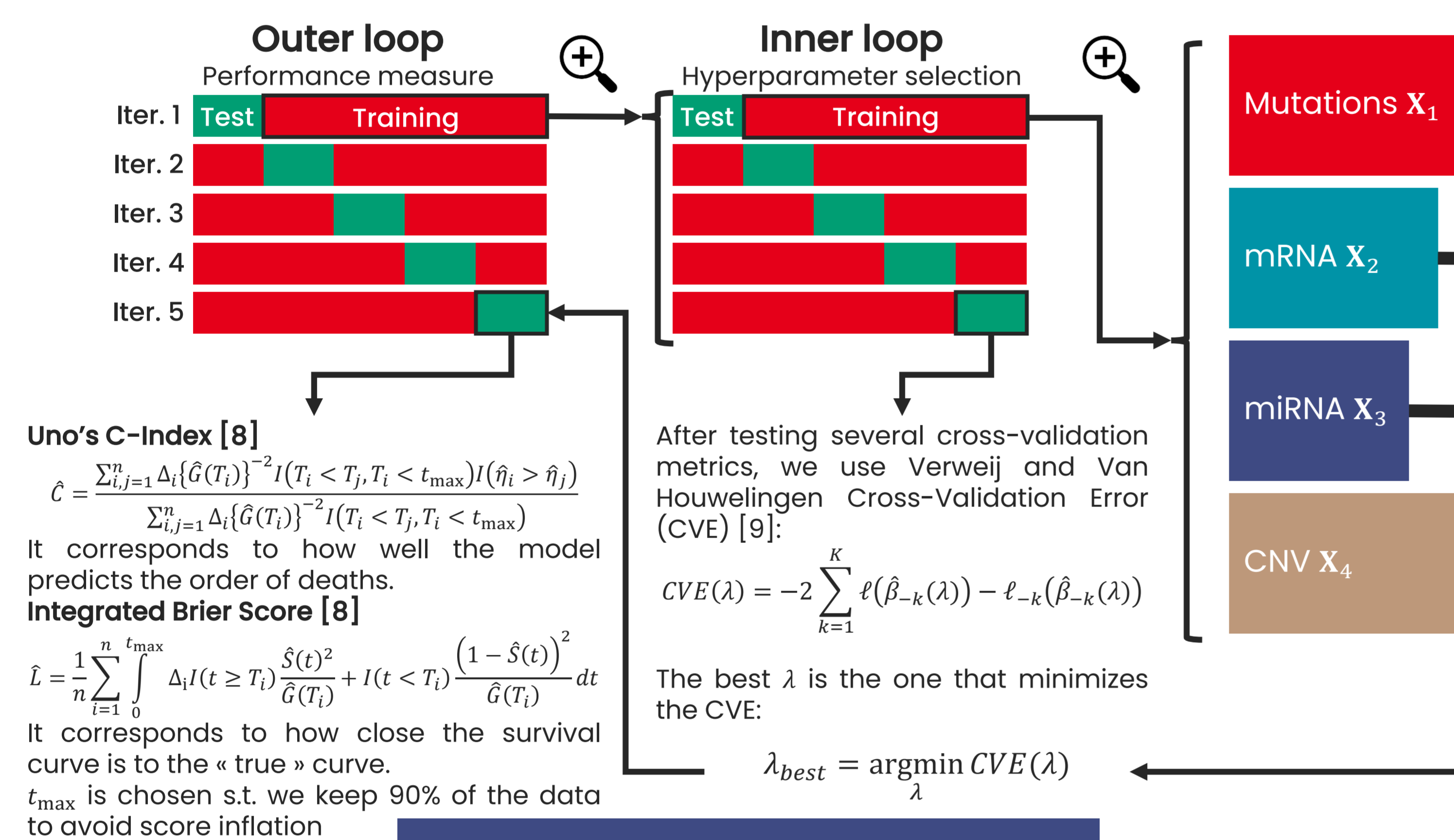
Supervised JIVE [7]

$$\operatorname{argmin}_{\mathbf{U}_1, \mathbf{S}_1, \mathbf{W}_1, \dots, \mathbf{U}_L, \mathbf{S}_L, \mathbf{W}_L} (1-\eta) \left\| \begin{bmatrix} \mathbf{X}_1^T \\ \vdots \\ \mathbf{X}_L^T \end{bmatrix} - \mathbf{U}_j - \begin{bmatrix} \mathbf{W}_1 \mathbf{S}_1 \\ \vdots \\ \mathbf{W}_L \mathbf{S}_L \end{bmatrix} \right\|_F^2 + \eta \left\| \mathbf{y}^T - \mathbf{S}_j^T \theta_1 - \sum_{l=1}^L \mathbf{S}_l^T \theta_{2l} \right\|_2^2$$

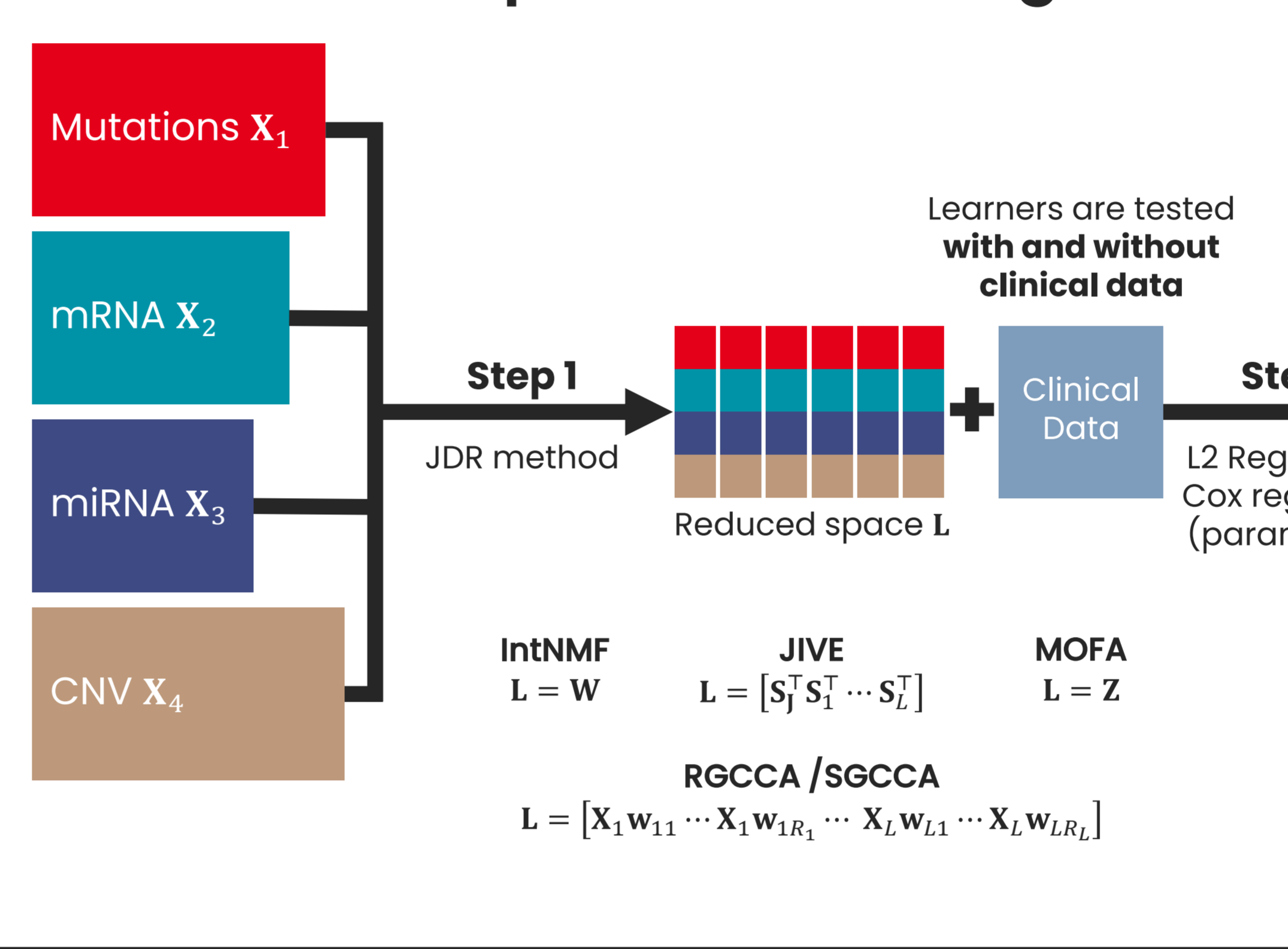
y: Deviance Residuals of Null Cox Model

Survival Prediction with JDR

Nested Cross-Validation



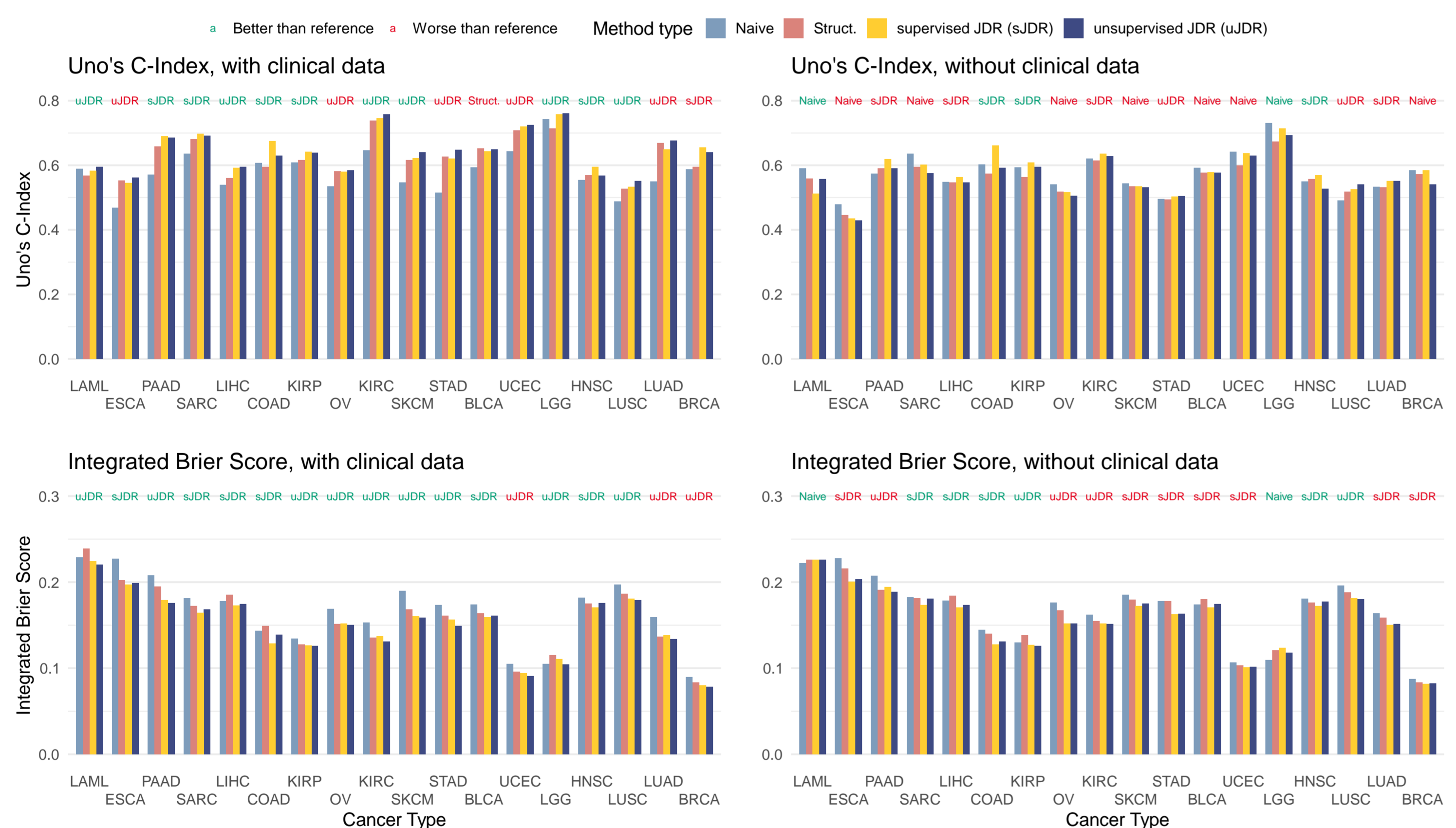
Sequential Modeling



Tested Parameters

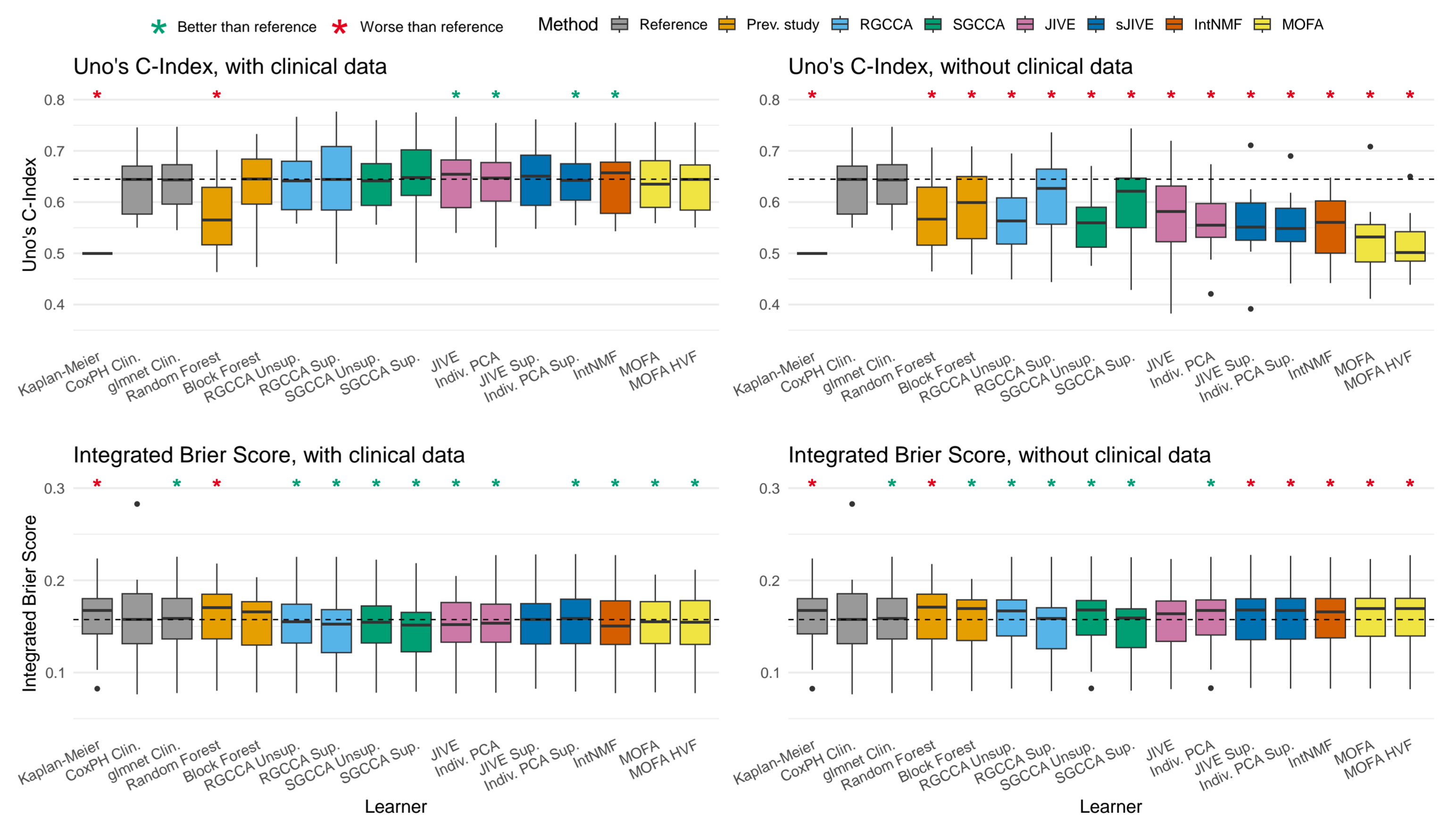
- RGCCA:**
 $\tau_l \in \{0, 1\}$; $R_l \in \{1, 3\}$; $g(x) = x^2$
 The 2 supervised/unsupervised connection graphs described above.
- SGCCA:**
 $s_l \in \{0, 5, 0, 75\}$; $R_l \in \{1, 3\}$; $g(x) = x^2$
 The 2 supervised/unsupervised connection graphs described above.
- JIVE:**
 $r_j \in \{0, 1, 3\}$ and $r_a \in \{0, 1, 3\}$ except $\{0, 0\}$
- sJIVE:**
 $r_j \in \{0, 1, 3\}$ and $r_a \in \{0, 1, 3\}$ except $\{0, 0\}$
 $\eta \in \{0, 1, 0, 5, 0, 9\}$
- IntNMF:**
 $R = \{2, 3, 5\}$
- MOFA:**
 $R = 3$
 With and without restriction to Highly Variable Features (HVF): we keep the top 10% most variable features of each block.
 With and without scaling blocks to unit variance.

Results by cancer



Barplots for each class of methods and each cancer type. All tested parameter combinations were averaged. The best class of method for each dataset is written above, and colored green if the average performance of the winning class is better than the mean performance of CoxPH Clin. We consider individual PCA (JIVE with $r_j = 0$) as structured and joint PCA (JIVE with $r_a = 0$) as a JDR method. Dataset are ordered by number of observations.

Results by method



Boxplots of each tested method. Several parameter combinations were tested for each JDR methods, only the best performing ones are represented here. P-values were computed using a paired Wilcoxon test between the folds of the considered methods and those of CoxPH Clin. Combinations kept: RGCCA Sup./Unsup.: $R_l = 1$, $\tau_l = 0, 1$, SGCCA Sup./Unsup.: $R_l = 1$, $s_l = 0, 5$, JIVE: $r_j = 3$, $r_a = 1$, PCA: $r_j = 0$, $r_a = 1$, JIVE Sup.: $r_j = 3$, $r_a = 1$, $\eta = 0, 9$ PCA Sup.: $r_j = 0$, $r_a = 1$, $\eta = 0, 9$ IntNMF: $R = 3$, MOFA: 3 factors, scaled MOFA HVF: 3 factors, scaled.

Conclusion

With clinical data

- JDR methods** (supervised and unsupervised) **improve both C-Index and IBS.**
- IntNMF has the best median performance** across datasets.

Without clinical data

- C-index and IBS are inferior across all methods.
- JDR are the best performing non reference methods.
- Supervised RGCCA and SGCCA** have the highest performance.
- Supervising JIVE does not seem to improve performance**, as is the case with RGCCA/SGCCA.

Perspectives

- Add **more JDR methods** in the benchmark (MCIA, iCluster, Scikit-Fusion).
- Use an **automatic procedure** to select the **best performing parameters** for each method.
- Include methods capable of extracting **Common, Local and Distinct (CLD)** components, in line with the good performance of JIVE.
- Design a new method** based on the best performing ones of the benchmark to exploit biological *a priori* information.

Dataset	Best method with Clinical Data		Best method without Clinical Data	
	C-Index	IBS	C-Index	IBS
BRCA	RGCCA Sup.	CoxPH Clin.	CoxPH Clin.	CoxPH Clin.
LUAD	sJIVE	CoxPH Clin.	CoxPH Clin.	CoxPH Clin.
LUSC	Joint PCA *	sJIVE	Joint PCA	Priority LASSO *
HNSC	RGCCA Sup. *	SGCCA Sup. *	RGCCA Sup. *	SGCCA Sup. *
LGG	LASSO *	MOFA HVF *	LASSO	MOFA HVF *
UCEC	CoxPH Clin.	CoxPH Clin.	CoxPH Clin.	CoxPH Clin.
BLCA	IntNMF *	SGCCA Sup. *	RGCCA Sup.	SGCCA Sup. *
STAD	sJIVE *	IntNMF *	Glimnet Clin. *	CoxPH Clin.
SKCM	JIVE *	sJIVE *	Glimnet Clin.	Glimnet Clin.
KIRC	Joint PCA *	Joint PCA *	Glimnet Clin.	CoxPH Clin.
OV	MOFA HVF	MOFA HVF	IntNMF	IntNMF
KIRP	IntNMF *	Block Forest	RGCCA Unsup.	Individual PCA
COAD	RGCCA Sup. *	RGCCA Sup. *	RGCCA Sup. *	RGCCA Sup. *
LIHC	Individual PCA *	Ranger *	RGCCA Sup.	Ranger *
SARC	RGCCA Sup. *	SGCCA Sup. *	LASSO	RGCCA Sup. *
PAAD	sJIVE *	SGCCA Unsup. *	Glimnet Clin. *	SGCCA Unsup. *
ESCA	SGCCA Unsup.	RGCCA Sup.	CoxPH Clin.	RGCCA Sup.
LAML	Block Forest *	MOFA HVF *	Ranger *	Ranger *

[1] Herrmann, M., Probst, P., Hornung, R., Jurinovic, V., & Boulesteix, A. L. (2021). Large-scale benchmark study of survival prediction methods using multi-omics data. *Briefings in bioinformatics*, 22(3), bbaa187. <https://doi.org/10.1093/bib/bbaa187>

[2] Tenenhaus, M., Tenenhaus, A., & Groenen, P. J. F. (2017). Regularized Generalized Canonical Correlation Analysis: A Framework for Sequential Multiblock Component Methods. *Psychometrika*, 10(1007), 9573-9573. Advance online publication. <https://doi.org/10.1007/s11336-017-9573-x>

[3] Lock, E. F., Hoadley, K. A., Marron, J. S., & Nobel, A. B. (2013). Joint and Individual Variation Explained (JIVE) for Integrated Analysis of Multiple Data Types. *The Annals of Applied Statistics*, 7(1), 523-542. <https://doi.org/10.1214/12-AOAS597>

[4] Chalise, P., & Fridley, B. L. (2017). Integrative clustering of multi-level omic data based on non-negative matrix factorization algorithm. *PLoS one*, 12(5), e0176278. <https://doi.org/10.1371/journal.pone.0176278>

[5] Argelaguet, R., Velten, B., Armi, D., Dietrich, S., Zenz, T., Marioni, J. C., Buettner, F., Huber, W., & Stegle, O. (2018). Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets. *Molecular systems biology*, 14(6), e8124. <https://doi.org/10.15252/msb.20178124>

[6] Bastien, P., Bertrand, F., Meyer, N., & Maumy-Bertrand, M. (2015). Deviance residuals-based sparse PLS and sparse kernel PLS regression for censored data. *Bioinformatics (Oxford, England)*, 31(3), 397-404. <https://doi.org/10.1093/bioinformatics/btu660>

[7] Palzer, E. F., Wendt, C. H., Bowler, R. P., Hersh, C. P., Sato, S. E., & Lock, E. F. (2022). sJIVE: Supervised Joint and Individual Variation Explained. *Computational statistics & data analysis*, 175, 107547. <https://doi.org/10.1016/j.csda.2022.107547>

[8] Rahman, M. S., Ambler, G., Choodari-Oskooei, B., & Omar, R. Z. (2017). Review and evaluation of performance measures for survival prediction models in external validation settings. *BMC medical research methodology*, 17(1), 60. <https://doi.org/10.1186/s12874-017-0336-2>

[9] Verweij, P. J., & Van Houwelingen, H. C. (1993). Cross-validation in survival analysis. *Statistics in medicine*, 12(24), 2305-2314. <https://doi.org/10.1002/sim.4780122407>