

Integrative Analysis of Multi-omics Data Improves Model Predictions: An Application to Lung Cancer

Erica Ponzi (✉ erica.ponzi@medisin.uio.no)

University of Oslo: Universitetet i Oslo <https://orcid.org/0000-0002-6089-3512>

Magne Thoresen

University of Oslo: Universitetet i Oslo

Therese Haugdahl Nøst

Universitetet i Tromsø: UiT Norges arktiske universitet

Kajsa Møllersen

Universitetet i Tromsø: UiT Norges arktiske universitet

Research article

Keywords: data integration, dimension reduction, joint and individual variance explained, multi-omics, prediction models

Posted Date: October 19th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-92731/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at BMC Bioinformatics on August 5th, 2021. See the published version at <https://doi.org/10.1186/s12859-021-04296-0>.

RESEARCH

Integrative analysis of multi-omics data improves model predictions: an application to lung cancer

Erica Ponzi^{1*}, Magne Thoresen¹, Therese Haugdahl Nøst² and Kajsa Møllersen²

Abstract

Background: Cancer genomic studies often include data collected from several omics platforms. Each omics data source contributes to the understanding of the underlying biological process via source specific ("individual") patterns of variability. At the same time, statistical associations and potential interactions among the different data sources can reveal signals from common biological processes that might not be identified by single source analyses. These common patterns of variability are referred to as "shared" or "joint". To capture both contributions of variance, integrative dimension reduction techniques are needed. Integrated PCA is a model based generalization of principal components analysis that separates shared and source specific variance by iteratively estimating covariance structures from a matrix normal distribution. Angle based JIVE is a matrix factorization method that decomposes joint and individual variation by permutation of row subspaces. We apply these techniques to identify joint and individual contributions of DNA methylation, miRNA and mRNA expression collected from blood samples in a lung cancer case control study nested within the Norwegian Woman and Cancer (NOWAC) cohort study.

Results: In this work, we show how an integrative analysis that preserves both components of variation is more appropriate than analyses considering uniquely individual or joint components. Our results show how both joint and individual components contribute to a better quality of model predictions, and facilitate the interpretation of the underlying biological processes.

Conclusions: When compared to a non integrative analysis of the three omics sources, integrative models that simultaneously include joint and individual components result in better prediction of cancer status and metastatic cancer at diagnosis.

Keywords: data integration; dimension reduction; joint and individual variance explained; multi-omics; prediction models

*Correspondence:

erica.ponzi@medisin.uio.no

¹Oslo Center for Biostatistics and Epidemiology, UiO, University of Oslo, Oslo, Norway

Full list of author information is available at the end of the article

1 Background

Cancer studies benefit from the availability of genomic data, also known as omics.

The dimensionality of omics data is extremely high, suggesting the application of

4 dimension reduction techniques. Additionally, omics are available across multiple
5 sources (or ‘blocks’) of data, collected on the same organisms or tissues, and mea-
6 sured on different platforms. A comprehensive understanding of the key underlying
7 biological process relies on an integrative approach able to combine the information
8 arising from such multi-source data. To this end, a large number of statistical meth-
9 ods for the simultaneous analysis of multi-omics data have recently been proposed.
10 Multiple reviews of such methods are available, for example in [1-3].

11 Data integration techniques are often used to identify ‘joint’ (also referred to as
12 ‘common’ or ‘shared’) contributions of the data sources to the observed variation,
13 and their simultaneous effect on the biological process under study. Such patterns
14 of variation arise from the interaction among different omics sources, and may not
15 be detected by a separate analysis of each single source. However, the different data
16 sources do not only contain the joint information, but also independent contribu-
17 tions. The separate analysis of each data source has so far been the most common
18 approach used in the omics context, and knowledge about the individual contribu-
19 tions of each omics source is relevant to the understanding of the biological processes
20 of interest. As a consequence, considering only the joint patterns might also prove
21 insufficient, as it overlooks the heterogeneity among single data sources, and their
22 individual signals from the underlying relevant biological process. An example of
23 this can be seen in genomic studies collecting DNA methylation and gene expres-
24 sion data. It is known that methylation regulates gene expression and that this
25 can cause a non-negligible joint structure across the different data sources. On the
26 other hand, methylation and gene expression correlate to the clinical outcome also
27 through signals that are specific to each omics data source and biologically relevant
28 independently from each other. Therefore, dimension reduction methods that take
29 both joint and individual patterns into account are necessary.

30 Two of the most common approaches to dimension reduction are Principal Com-
31 ponent Analysis (PCA) and matrix factorization, and these have recently been
32 expanded to the case of multi source data. For example, consensus PCA [4] consists
33 of PCA on the normalized concatenated data, and distributed PCA [5] performs
34 local PCA on the individual data sources and then uses these principal components
35 to estimate a global covariance structure. Other dimension reduction methods have
36 been extended to the case of multi source data, as for example canonical correlation

37 analysis (CCA) [6] or partial least squares (PLS) analysis, which has been further
38 generalized to O2PLS [7]. A similar method that allows for the presence of multiple
39 data sources is the multiple CCA [8], but it mainly focuses on the common varia-
40 tion among the components, and seems to neglect the individual contributions of
41 the data sources.

42 Although the extensions of PCA mentioned above account for multiple data
43 sources, they focus either on their joint or individual contributions but do not
44 identify both of them simultaneously. Integrated PCA (iPCA) is a model based
45 generalization of PCA that decomposes variance into shared and source specific
46 variation [9]. It is based on the assumption of a matrix variate normal distribu-
47 tion of the data, whose covariance structure is given as the Kronecker product of
48 a shared and an individual, data block specific, covariance matrix. Integrated PCA
49 estimates these two matrices via an iterative algorithm, and consequently extracts
50 principal components for both.

51 Solutions based on matrix factorization also preserve both joint and individual
52 structures in the data. In this framework, each data block is decomposed into three
53 matrices modeling different types of variation, specifically joint variation across the
54 blocks, individual variation for each data block, and residual noise. JIVE stands
55 for Joint and Individual Variance Explained, it was formulated by [10] and, also
56 thanks to its implementation available in R [11], has been used in various medical
57 applications, including clustering of cancer genomics data [12], multisource omics
58 data [13, 14] and imaging and behavioral data [15]. Although JIVE solves the issue
59 of maintaining joint and individual structures, it uses an iterative algorithm and is
60 computationally very intensive. In [16], Angle Based JIVE (aJIVE) was formulated
61 to improve this aspect. It computes the matrix decomposition by using perturba-
62 tions of the row spaces to identify the joint and individual variation, and results
63 in a much faster implementation than the original JIVE. Other similar approaches
64 to identify both kinds of variation have been proposed, as for example DISCO [17]
65 and OnPLS [18]. An illustration of these methods and a comparison with JIVE was
66 provided in [19].

67 In this work, we demonstrate how an integrative analysis of the data is more
68 appropriate than a uniquely source-specific analysis, which would not reveal rela-
69 tions between the sources of biological information. Integrative analyses are also

70 preferable to uniquely joint analysis, which can be redundant and use the available
71 information to repeatedly measure the same shared components. We show how both
72 joint and individual components contribute to a better quality of model predictions.
73 The combination of joint and individual components can also facilitate the biolog-
74 ical interpretation of the underlying process, although this might still fail as the
75 dimension reduction itself bears the risk of obscuring some relevant information.

76 We use both a matrix factorization and a PCA based method to identify indi-
77 vidual and joint components in a real data set on lung cancer. We chose to use
78 aJIVE among the matrix factorization methods, because it inherits a good sub-
79 space recovery in comparison to other methods [19], as well as the robustness to
80 model misspecification from JIVE, but it also solves the issue of correlated indi-
81 vidual subspaces [16], and provides a much faster implementation. Furthermore,
82 [20] show that aJIVE performs best in terms of consistency and lack of overfitting
83 when compared to other integrative methods. We validate the aJIVE results with a
84 PCA based method, which has the advantage of being independent from the initial
85 ranks selection. We use iPCA because it is the only method in this framework that
86 can identify both individual and joint patterns, and because it has been shown to
87 perform well in prediction models [9].

88 The data we use stem from a lung cancer case control study nested within the Nor-
89 wegian Woman and Cancer (NOWAC) cohort study [21]. The associations among
90 three levels of omics data analyzed in blood samples, specifically DNA methyla-
91 tion, mRNA and miRNA expression, are investigated and their joint and individual
92 contributions are used to predict future cancer cases, and for the characterization
93 of future cancers as metastatic or non-metastatic at diagnosis. We show that both
94 kinds of components contain information that reveal properties about biological
95 processes and that using joint components improves model predictions. Specifically,
96 we show that joint components increase the AUC of prediction models, when com-
97 pared to models with solely individual components, to models uniquely based on
98 clinical, patient-level covariates, and most importantly to non-integrative models,
99 i. e. based on independent analyses of data from each source.

100 **2 Methods**

101 2.1 Data integration setup

102 Throughout the manuscript, we will denote each data block with \mathbf{X}_k , where $k =$
103 $1, \dots, K$ and K is the number of data sources used in the study. Each block is a matrix
104 with n columns, where n is the number of study subjects. The k th matrix \mathbf{X}_k has
105 p_k rows, corresponding to the variables in data source k . The overall dimensionality
106 is denoted as $p = p_1 + \dots + p_K$. The low-rank decomposition we want to obtain from
107 the different methods is:

$$\begin{aligned} \mathbf{X}_1 &= \mathbf{J}_1 + \mathbf{I}_1 + \boldsymbol{\epsilon}_1 \\ &\vdots \\ \mathbf{X}_K &= \mathbf{J}_K + \mathbf{I}_K + \boldsymbol{\epsilon}_K \end{aligned} \tag{1}$$

where \mathbf{I}_k is the individual component for data block k , $\boldsymbol{\epsilon}_k$ is its residual component
and

$$\mathbf{J} = \begin{bmatrix} \mathbf{J}_1 \\ \dots \\ \mathbf{J}_K \end{bmatrix} \tag{2}$$

108 is the joint structure matrix, where each \mathbf{J}_k is the submatrix of the joint structure
109 \mathbf{J} associated with \mathbf{X}_k . The methods we use in our application obtain such decom-
110 position with very different algorithms, but both focus on the distinction between
111 the \mathbf{J} and \mathbf{I}_k components, and on the calculation of their loadings and scores.

112 2.2 Angle Based JIVE

113 Angle based Joint and Individual Variation Explained (aJIVE) is a variant of the
114 JIVE method, based on perturbation of row subspaces. JIVE aims to minimize
115 the squared residual components $\epsilon_1, \dots, \epsilon_K$, using an iterative algorithm that al-
116 ternatively estimates the joint and individual components by singular value de-
117 composition (SVD). AJIVE builds on this method but constructs the algorithm
118 in a more efficient and computationally feasible way. Besides resulting in a faster
119 implementation of the algorithm, aJIVE provides a more intuitive interpretation

120 of the decomposition, especially in the case of high correlations among individual
 121 components [16]. The aJIVE algorithm is structured in three phases: First the low-
 122 rank approximation of each data block \mathbf{X}_k is obtained by SVD. Secondly, the joint
 123 structure between the obtained low-rank approximations is extracted computing
 124 the SVD of the stacked row basis matrices, by using basic principles of Principal
 125 Angle Analysis. Finally, the joint components \mathbf{J}_k are obtained by projection of each
 126 data block onto the joint basis, while the individual components \mathbf{I}_k are calculated
 127 by orthonormal basis subtraction.

128 The first step is based on the choice of the initial ranks, which are used as a
 129 threshold value in the first SVD decomposition of the data blocks. This choice is
 130 rather subjective and involves taking into account some bias variance trade-off in the
 131 joint signals representation. Although [16] provide guidelines on how to determine
 132 the initial ranks, the recommended choice is based on the observation of scree plots,
 133 which remains highly subjective. As an alternative, [22] present a choice of initial
 134 ranks based on the profile likelihood of the single data blocks.

135 From the aJIVE decomposition, it is possible to obtain the full matrix represen-
 136 tation of the original features, as well as the block specific decompositions of each
 137 data source and the common normalized scores. The algorithm requires much lower
 138 computation time than the original JIVE, and its implementation is available in
 139 Matlab [23] and R [24].

140 2.3 Integrated PCA

141 Integrated Principal Component Analysis (iPCA) is a generalization of principal
 142 component analysis to multiple data sources. The core idea of iPCA is the assump-
 143 tion that each data source \mathbf{X}_k follows a matrix variate normal distribution:

$$\mathbf{X}_k \sim \mathbf{N}_{(n,p)}(\mathbf{M}, \mathbf{\Sigma}, \mathbf{\Delta}_k), \quad (3)$$

144 where \mathbf{M} is the mean matrix, $\mathbf{\Delta}_k$ is the column covariance specific to \mathbf{X}_k and $\mathbf{\Sigma}$
 145 is the row covariance shared by all data matrices \mathbf{X}_k . The assumption of normality
 146 is very common for methylation data, and for gene expression after applying a log
 147 transformation. Moreover, [9] show that the iPCA method is robust to deviations
 148 from normality.

Integrated PCA performs PCA on the covariance matrices Σ and Δ_k and denotes with $\hat{\mathbf{U}}$ the eigenvectors of Σ , which represent the joint components \mathbf{J} . The eigenvectors of Δ_k are denoted as $\hat{\mathbf{V}}_k$ and represent the individual components \mathbf{I}_k specific to data set k . To obtain these quantities, iPCA uses the “Flip-Flop algorithm” to calculate maximum likelihood estimators (MLEs) of the covariance matrices. Because of the complexity of the model and of the non-existence of MLEs for a large number of cases, iPCA uses a penalized maximum likelihood, with a penalty term defined as:

$$P_q(\Sigma^{-1}, \Delta_1^{-1}, \dots, \Delta_k^{-1}) = \lambda_\Sigma \|\Sigma^{-1}\|_q + \sum_{k=1}^K \lambda_k \|\Delta_k^{-1}\|_q$$

149 with norm q . Different possibilities for q are illustrated in [9]: L_1 or multiplicative
 150 Frobenius or additive Frobenius, and suggest using the multiplicative Frobenius
 151 to ensure convergence. The choice of the tuning parameters λ_Σ and $\lambda_1, \dots, \lambda_k$
 152 is based on a missing data imputation framework. This step of the algorithm is
 153 computationally very intensive and requires the input of initial values, which need
 154 to be chosen carefully by the user.

155 Integrated PCA has been implemented in R and its code is available at [25]. Note
 156 that iPCA, unlike aJIVE, does not require the specification of the initial ranks, al-
 157 though it requires initial values for the tuning parameters. While aJIVE is built on
 158 an optimization procedure, focused on minimizing the approximation error, iPCA is
 159 a model-based approach for the estimation of the underlying subspace. This ensures
 160 the convergence to a global solution for iPCA, and a better recovery of the subspace,
 161 but on the other hand it results in less robustness to model misspecification, which
 162 does not affect aJIVE, and less robustness to errors. Both aJIVE and iPCA do not
 163 rely on any independence assumptions, and correlations among individual compo-
 164 nents do not represent an issue in any of the two methods. An important difference
 165 between the two methods, as we highlighted above, is the computation time of the
 166 two algorithms, which is significantly lower in aJIVE. The interpretability of the
 167 results is comparable in aJIVE and iPCA.

168

169 2.4 Application to the NOWAC data

170 *The dataset*

171 The dataset used in the following analyses stems from blood samples in a lung
172 cancer case control study nested within the Norwegian Women and Cancer Study
173 (NOWAC) [21]. All participating subjects are women who did not have a cancer di-
174 agnosis at recruitment (1991–2006) or at time of blood sampling (2003–2006). The
175 time from blood sampling to cancer diagnosis ranges from 0.3 to 7.9 years, with a
176 median time equal to 4.2 years. For each case, one control was matched on time
177 since blood sampling and birth year. All participants gave written informed consent
178 and the study was approved by the Regional Committee for Medical and Health
179 Research Ethics and the Norwegian Data Inspectorate. Three levels of omics data
180 are available for $n = 230$ individuals, with numbers of variables respectively equal
181 to $p_1 = 485512$ CpG methylation, $p_2 = 11610$ mRNA expression and $p_3 = 199$
182 miRNA expression. Information about individual covariates, including age, body
183 mass index (BMI), dietary and smoking habits was also collected for all partici-
184 pants. Outcomes of interest are the classification of case versus control, as well as
185 the characterization of cancers as metastatic or non-metastatic at diagnosis.
186 Both methylation and gene expression have been shown to associate with the oc-
187 currence and characteristics of lung cancer [26–29] but it is also known that the
188 different data sources might contribute together and jointly relate to these bio-
189 logical outcomes [30, 31]. We apply aJIVE and iPCA to estimate such joint and
190 individual contributions, and use these components in prediction models for the
191 occurrence of lung cancer and classification of tumor types.

192 *Filtering and preprocessing*

193 Laboratory processing and microarray analyses for mRNA expression and DNA
194 methylation are described in [31]. For miRNA, laboratory processing included
195 miRNA isolation and purification from 100 μ l plasma using the Qiagen miRNeasy
196 Serum/Plasma Kit. Small RNA sequencing libraries were prepared using the
197 NEXTflex small RNA-seq kit v3 (Bioo Scientific, Austin, TX, USA) and sequencing
198 of fragments was performed using a Illumina HiSeq4000 flowcell, according to the
199 manufacturer’s instructions (Illumina, Inc., San Diego, CA, USA), at 50 bp SE,
200 resulting in approximately 7 – 9M reads per sample.

201 The filtering of miRNA expressions was based on the counts per million, that is
202 the total read counts of a miRNA divided by the total read counts of the sample and
203 multiplied by 10^6 , and signals having less than one count per million were excluded.
204 Additionally, signals with null reads on more than 5 patients were excluded.

205 Because of the high computational requirements, we reduced the number of
206 mRNA expressions to $p_2 = 5000$, by selecting the variables with higher variance. We
207 then reduced the number of CpGs methylation sites by selecting the CpGs located
208 on the same genes as the filtered mRNAs. Among these, we excluded CpGs with
209 more than 40% missing data, as well as CpGs with extreme M-values ($|M| > 3$,
210 see [32, 33]). This resulted in $p_1 = 18545$. All $p_3 = 199$ available miRNAs were
211 analysed. We used log2 transformed expressions for both mRNA and miRNA, and
212 M-values for methylation [34]. We accounted for missing values in the data by using
213 `SVDmiss`, as suggested in [10], and centered the data for easier interpretation and
214 comparison of the two methods.

215 *aJIVE and iPCA*

216 We performed aJIVE on the three levels of omics data. The initial ranks were
217 selected by maximizing the profile likelihood [22], and set respectively to 43, 11
218 and 9. Different choices of initial ranks were also explored. Additionally, we ran
219 aJIVE on pairs of omics sources, with the same ranks as above for each two way
220 association.

221 We performed iPCA on a smaller subset of variables, due to the high compu-
222 tational requirements. We selected variables with the same procedure as reported
223 above, but starting with $p_2 = 500$ mRNAs, and resulting in $p_1 = 1588$ CpGs. The
224 $p_3 = 199$ miRNAs were all included. We also repeated aJIVE on this smaller subset.
225 We set the initial penalty terms to 1, 0.01 and 0.001, as suggested in [9], and used
226 the multiplicative Frobenius norm.

227 *Using joint components for prediction*

228 Joint and individual components were used in prediction models. The outcome of
229 interests were the occurrence of lung cancer (yes/no) and metastasis (yes/no).

230 We fitted logistic models on each outcome using joint and individual components
231 as explanatory variables, in addition to age, BMI and smoking. These models were
232 compared in terms of AUC with the respective models with only age, BMI and

233 smoking as covariates. To assess the performance of the models, we measured the
234 average AUC in a 10 fold cross validation. To identify the amount of information
235 added by the joint components, we compared these models to models fitted on the
236 individual components from aJIVE and the clinical covariates as predictors. Finally,
237 we compared these results with a non-integrative analysis, obtained by performing
238 PCA separately on each single data source. We fitted a model on the first principal
239 components (PCs) of each data source, and on the same clinical covariates. We chose
240 to include five PCs for each data source, to have a comparison with the number of
241 joint and individual components, and based on the variance explained by the first
242 PCs and on the analysis of screeplots. Models were repeated with a higher number
243 of PCs and results did not change substantially.

244 In addition, we used a random forest of 500 trees to predict case vs control on
245 the basis of joint and individual components, and patient covariates as above (age,
246 BMI, smoking). We extracted the AUCs and the out of the bag (OOB) classification
247 errors from the random forest, and we ranked all the variables in importance on the
248 basis of their mean decrease in gini index [35].

249 **3 Results**

250 We will focus on aJIVE as this is where we have the lowest computational re-
251 quirements and can include a higher number of variables from each data source. A
252 comparison with iPCA will be given in Section 3.3.

253 **3.1 aJIVE**

254 Using initial ranks obtained with the profile likelihood method resulted in a joint
255 rank equal to 5, and individual ranks respectively equal to 38, 7 and 7. The de-
256 composition is illustrated in the heatmap in Figure 1 and shows that the joint
257 component dominates the individual and residual components for all three data
258 sources. Figure 2 reports the proportions of variance explained that are due to the
259 joint, individual and residual components and also shows that the joint components
260 explain most of the total variance in the data.

261 Repeating the algorithm with lower initial ranks resulted in lower joint rank esti-
262 mates, while higher ranks gave estimates of joint ranks that were stable and equal
263 to 5, suggesting that overestimating the initial ranks should not affect the results,
264 while an underestimation of the ranks can fail to detect some joint components

265 in the data. The pairwise analysis starting with the same initial ranks resulted
266 in different joint ranks, suggesting that most joint components are shared between
267 methylation and mRNA. The aJIVE joint and individual ranks for pairwise analyses
268 are reported in Table 1.

269 3.2 Prediction models

270 Figure 3 reports the ROC curves relative to the logistic models fitted on the full
271 dataset. The model with only patient covariates (age, BMI and smoking) as ex-
272 planatory variables is compared in terms of AUC to three "integrative" models,
273 reported in Table 2: the model with aJIVE joint components and patient covari-
274 ates as explanatory variables, the model with aJIVE individual components and
275 patient covariates, and finally the "full" model with patient covariates, aJIVE joint
276 components and the first five aJIVE individual components for each data source
277 as explanatory variables. These are further compared to "non-integrative" models,
278 using the first five individual PCs obtained for each dataset separately.

279 These results were validated by 10 fold cross validation for each outcome. In the
280 ROC studies from cross validation, the model with all components seems to improve
281 the prediction for both case-control and metastasis status. The mean AUCs for the
282 "full" models are 0.76 and 0.74, for case-control and metastasis status respectively.
283 The mean AUCs for the models using only joint components and clinical covariates
284 are 0.76 and 0.71, while the mean AUCs for the models using only clinical covariates
285 are 0.69 and 0.69. We see that including the individual components does not always
286 result in better models, which seems to confirm that the major role is played by the
287 joint components identified by aJIVE rather than the source specific components.
288 When including only individual components from aJIVE, in addition to clinical
289 covariates, the mean AUC is 0.66 for the prediction of case vs control and 0.68 for
290 the prediction of metastasis. The mean AUC of the non-integrative model, based
291 on the single data PCAs and the clinical covariates, is respectively 0.68 and 0.65,
292 lower than the AUCs obtained in the models using the aJIVE components.

293 Table 2 reports accuracy and OOB classification error for the random forests,
294 as well as the mean AUCs. For case-control status, both models with all compo-
295 nents and with only joint components improve the predictions. The non-integrative
296 models perform similarly to the model including only covariates, and have very low

297 accuracy. This difference from the logistic models with cross validation can be due
298 to the instability of the random forest, and to the limited sample size. We do not
299 report the random forests results for metastasis because they are highly unstable
300 and the accuracy is very low, most likely due to the even more limited sample size,
301 that is only 125 (only cases) for the metastasis.

302 Figure 4 shows the first ten variables ranked by variable importance in the full
303 model for case-control status. Three of the five joint components appear among
304 the first five variables when ranked for variable importance in the random forest
305 prediction.

306 3.3 iPCA

307 Joint components obtained from iPCA were used in prediction models of case-
308 control and metastasis status. Joint components were also obtained by aJIVE on
309 the same subset of data and prediction models were compared. The initial ranks
310 were selected for aJIVE via profile likelihood and set respectively to 23, 18 and 9.
311 AJIVE selected a joint rank equal to 6 and individual ranks equal to 18, 13 and 6.

312 We fitted prediction models on case-control and metastasis status using the joint
313 scores from iPCA and aJIVE as predictors, together with the patient covariates.
314 Figure 5 shows ROC curves for both methods for each outcome. In both cases, using
315 joint scores from either aJIVE or iPCA improved the prediction of the outcome to
316 an equivalent extent. In the prediction of case-control status, aJIVE resulted in
317 slightly higher AUCs than those based on iPCA, both when using five and 10 joint
318 components from iPCA. When predicting metastasis, aJIVE performed better than
319 iPCA with five joint components, but not with 10 joint components. However, the
320 results were very similar. Results were validated via 10-fold cross validation and the
321 AUC are comparable between iPCA and aJIVE based models. The mean AUC for
322 the models including five iPCA components and patient covariates is 0.72 for case-
323 control status and 0.70 for metastasis status, while including 10 iPCA components
324 gives mean AUC equal to 0.71 and 0.73 respectively. The mean AUC of the model
325 including aJIVE components and patient covariates is 0.76 for case-control status
326 and 0.73 for metastasis status.

327 4 Discussion

328 We use two data integration methods to identify both joint and individual com-
329 ponents in a lung cancer study, where multiple omics data sources are available.
330 While the individual contribution of each data source is known to be relevant and
331 has been widely studied in this context, different data sources are also expected to
332 jointly associate with the clinical outcomes. We show how including both joint and
333 individual components in prediction models improves the quality of prediction for
334 the occurrence of lung cancer, as well as its classification into metastatic or non-
335 metastatic cancer. Models that include both components lead to better predictions
336 when compared to non-integrative models, or to models based on clinical covariates.

337 Prediction models are validated in a 10 fold cross validation framework, and such
338 results are further confirmed by random forests. From the cross validation study,
339 we see that for case-control status, the joint components provide better prediction
340 than non-integrative analysis. For metastasis, the improvement is evident only for
341 the "full" integrative model with both joint and individual components. Still, the
342 AUCs are comparable and small differences can be due to the limited sample size,
343 since we are restricted to cases only ($n = 125$), and to the unbalance between groups
344 in the analysis of metastasis. This limitation is more obvious in the random forests,
345 where the small sample size leads to low accuracy and high instability. Therefore,
346 we only show random forests for the classification of cases vs controls, although the
347 same procedure can be repeated for the prediction of metastasis.

348 A possible explanation of generally low AUCs is that prediction models might also
349 be affected by the time between blood sampling and cancer diagnosis, and we expect
350 the quality of predictions to be higher in subjects with a shorter time to diagnosis.
351 We stratified cases into two subgroups based on time to diagnosis (higher vs lower
352 than the median time) and obtained higher in-sample AUCs for the classification
353 of case vs control in subjects with a closer time to diagnosis for most models. For
354 the classification of metastasis, the sample size in the two time to diagnosis classes
355 is not enough to draw conclusions.

356 It is interesting to observe that three genomic components identified from aJIVE
357 rank above smoking in importance for case-control classification (Figure 4). Smok-
358 ing is known to be the one, major risk factor for lung cancer. We speculate that
359 one reason might be a weak mediation of the smoking effect through the genomic

360 components [36]. While this work provides a preliminary evidence of the impor-
361 tance of an integrative analysis of the omics sources, a more thorough investigation
362 of the joint and individual components could help identify relevant biological pat-
363 terns for future research. An example can be given by the underlying biological
364 processes involving smoking and lung cancer: the omics signals that are dominat-
365 ing the components could be important risk factors for lung cancer, in addition to
366 information on active or past smoking, and their interaction could shed light on
367 the relevant underlying biological processes. Although a functional interpretation
368 of such processes and of their link to the clinical outcomes is not straightforward,
369 an investigation of the aJIVE components could provide further information that
370 would not be identified by a non integrative analysis of the separate omics sources.

371 The chosen approach for variable filtering is based on variance for mRNA, and on
372 genomic location for methylation. Specifically, the top 5000 most variable mRNAs
373 are selected and CpGs are then selected based on their location on genes, by includ-
374 ing CpGs located on the same genes as filtered mRNAs. The joint contributions are
375 therefore expected to be most relevant, given the natural association between sig-
376 nals on the same gene location, and this aspect needs to be carefully considered in
377 the interpretation of the results. Also the filtering of the miRNAs needs to be taken
378 into account, where less restrictive criteria might result in the estimation of differ-
379 ent joint and individual components. Other choices could be made in this phase, for
380 example applying the variance criterion independently on each data source, which
381 could yield different joint and individual components. Alternative criteria are the
382 interquartile range (IQR), or the association with the clinical outcome of interest,
383 estimated by an appropriate regression model. Another choice we made in the pre-
384 processing and filtering of the data is the use of M-values for methylation. This
385 choice is motivated by [34], but β values could also be employed for the purpose of
386 this work.

387 To identify joint and individual components, we use one matrix factorization based
388 method, aJIVE, and one PCA based method, iPCA. Although iPCA requires higher
389 computation time when compared to aJIVE, it solves one of the main issues in
390 aJIVE, which is the selection of initial ranks. The most common method for the
391 choice of initial ranks in aJIVE is the visualization of screeplots, which is subjective
392 and highly sensitive to noise in the data. The profile likelihood idea suggested by

393 [22] partly addresses the problem, but it still lacks some objectivity and automa-
394 tion. Nevertheless, the correct choice of ranks is fundamental for aJIVE, and ranks
395 misspecification can lead to incorrect results [16]. For this reasons, iPCA should
396 be preferred. On the other hand, the implementation and computation of aJIVE
397 is much more efficient and fast, and allows for inclusion of more variables in the
398 analysis. Our results also show that aJIVE gives slightly better predictions than
399 iPCA. Therefore, a sensible strategy could be to use aJIVE when the choice of ini-
400 tial ranks is non-controversial and quite straightforward from the data, and to opt
401 for iPCA when such choice is uncertain. In any case, it is always recommended to
402 fit the models with more than one method and compare and validate results.

403 The high dimensionality of the data also motivates the use of sparse methods,
404 which reduce the number of variables included in the model and provide an easier
405 interpretation of the results. A sparse version of the aJIVE method could be used
406 for this purpose, by introducing a penalty term in the decomposition to induce
407 variable sparsity. This has not been specifically implemented for aJIVE, but [10]
408 discuss and provide an implementation of a sparse version of the JIVE method. It
409 has been shown in [9] that sparsity can be achieved in the iPCA method by using
410 an additive L1 norm in the estimation of the covariance matrices.

411 One aspect that is not accounted for by either methods is the presence of partially
412 shared components. When joint components are only shared by, for example, two
413 out of the three data sources, they will not be identified by aJIVE or iPCA. This
414 is a limitation of most data integration methods, and we expect partially shared
415 components to result in even better prediction models. A way to investigate par-
416 tially shared patterns is provided in the SLIDE method by [37], and is a potential
417 starting point for further work in this direction.

418

419 **5 Conclusion**

420 Our study shows how integrative models that include both joint and individual
421 contribution of multiple datasets lead to more accurate model predictions, and
422 facilitate the interpretation of the underlying biological processes. We use iPCA and
423 aJIVE to identify joint and individual contributions of DNA methylation, miRNA
424 and mRNA expression in a lung cancer case control study. We include both joint

425 and individual components in prediction models for the occurrence of lung cancer,
426 and for its classification into metastatic or non-metastatic cancer. Our results show
427 that models that include joint and individual components lead to better predictions
428 when compared to non-integrative models, or to models based on clinical covariates
429 only.

430 **Ethics approval and consent to participate**

431 All participants gave written informed consent and the study was approved by the Regional Committee for Medical
432 and Health Research Ethics and the Norwegian Data Inspectorate. More information is available in [21].

433 **Consent for publication**

434 Not applicable.

435 **Availability of data and materials**

436 The code for the statistical analysis is available at <https://github.com/ericaponzi>.

437 Data cannot be shared publicly because of local and national ethical and security policies. Data access for
438 researchers will be conditional on adherence to both the data access procedures of the Norwegian Women and
439 Cancer Cohort and the UiT –The Arctic University of Norway (contact via Tonje Braaten tonje.braaten@uit.no and
440 Arne Bastian Wiik arne.b.wiik@uit.no) in addition to an approval from the local ethical committee.

441 **Competing interests**

442 The authors declare that they have no competing interests.

443 **Funding**

444 Norwegian Research Council - grant number 248804: National training initiative to make better use of biobanks and
445 health registry data.

446 Norwegian Research Council - FRIMEDBIO grant number 262111: Identifying biomarkers of metastatic lung cancer
447 using gene expression, DNA methylation and microRNAs in blood prior to clinical diagnosis (Id-Lung).

448 **Acknowledgements**

449 The miRNA and mRNA analyses were provided by the Genomics Core Facility (GCF), Norwegian University of
450 Science and Technology (NTNU). GCF is funded by the Faculty of Medicine and Health Sciences at NTNU and
451 Central Norway Regional Health Authority.

452 **Author's contributions**

453 K. M., M. T. and E. P. conceived the research idea. E.P. conducted the statistical analyses. T. H. N. was
454 responsible for the acquisition of data and the biological interpretation of the results. E.P. K. M. and M.T. wrote
455 the manuscript, with inputs from T. H. N. All authors gave final approval.

456 **Author details**

457 ¹Oslo Center for Biostatistics and Epidemiology, UiO, University of Oslo, Oslo, Norway. ²Department of
458 Community Medicine, UiT, The Arctic University of Norway, Tromsø, Norway.

459 **References**

- 460 1. Tseng, G., Ghosh, D., Zhou, X.J.: Integrating Omics Data. Cambridge University Press, Cambridge (2015)
- 461 2. Huang, S., Chaudhary, K., Garmire, L.X.: More is better: Recent progress in multi-omics data integration
462 methods. *Frontiers in Genetics* **8** (2017)
- 463 3. Rappaport, N., Ron, S.: Multi-omic and multi-view clustering algorithms: review and cancer benchmark.
464 *Nucleic acids research* **42**, 10546–10562 (2018)
- 465 4. Westerhuis, J.A., Kourti, T., MacGregor, J.F.: Analysis of multiblock and hierarchical PCA and PLS models.
466 *Journal of Chemometrics* **12**, 301–321 (1998)

467 5. Fan, J., Wang, D., Wang, K., Zhu, Z.: Distributed estimation of principal eigenspaces. *Annals of Statistics* **47**,
468 3009–3031 (2019)

469 6. Hotelling, H.: Relations between two sets of variates. *Biometrika* **28**, 321–377 (1936)

470 7. Trygg, J., Wold, H.: O2-PLS, a two-block (x - y) latent variable regression (LVR) method with an integral OSC
471 filter. *Journal of Chemometrics* **17**, 53–64 (2003)

472 8. Witten, D., Tibshirani, R.: Extensions of sparse canonical correlation analysis with applications to genomic
473 data. *Statistical Applications in Genetics and Molecular Biology* **8**, 28 (2009)

474 9. Tang, T.M., Allen, G.I.: Integrated Principal Components Analysis (2018). [1810.00832](https://doi.org/10.100832)

475 10. Lock, E.F., Hoadley, K.A., Marron, J.S., Nobel, A.B.: Joint and individual variation explained (JIVE) for
476 integrated analysis of multiple data types. *Annals of Applied Statistics* **7**, 523–542 (2013)

477 11. O’Connell, M.J., Lock, E.F.: R.JIVE for exploration of multi-source molecular data. *Bioinformatics* **32**(18),
478 2877–2879 (2016). doi:[10.1093/bioinformatics/btw324](https://doi.org/10.1093/bioinformatics/btw324).
479 <https://academic.oup.com/bioinformatics/article-pdf/32/18/2877/25416720/btw324.pdf>

480 12. Hellton, K.H., Thoresen, M.: Integrative clustering of high-dimensional data with joint and individual clusters.
481 *Biostatistics* **17**(3), 537–548 (2016). doi:[10.1093/biostatistics/kxw005](https://doi.org/10.1093/biostatistics/kxw005).
482 <https://academic.oup.com/biostatistics/article-pdf/17/3/537/17741839/kxw005.pdf>

483 13. Kuligowski, J., Perez-Guaita, D., Sanchez-Illana, A., Leon-Gonzalez, Z., de la Guardia, M., Vento, M., Lock,
484 E.F., Quintas, G.: Analysis of multi-source metabolomic data using joint and individual variation explained
485 (JIVE). *Analyst* **13**, 4521–4529 (2015)

486 14. Kaplan, A., Lock, E.F.: Prediction with dimension reduction of multiple molecular data sources for patient
487 survival. *Cancer Inform* **16**, 1–11 (2017)

488 15. Yu, Q., Risk, B.B., Zhang, K., Marron, J.S.: JIVE integration of imaging and behavioral data. *NeuroImage* **152**,
489 38–49 (2017). doi:[10.1016/j.neuroimage.2017.02.072](https://doi.org/10.1016/j.neuroimage.2017.02.072)

490 16. Feng, Q., Jiang, M., Hannig, J., Marron, J.S.: Angle-based joint and individual variation explained. *Journal of*
491 *Multivariate Analysis* **166**, 241–265 (2018). doi:[10.1016/j.jmva.2018.03.008](https://doi.org/10.1016/j.jmva.2018.03.008)

492 17. Schouteden, M., Van Deun, T.F., Wilderjans, T.F., Van Mechelen, I.: Performing DISCO-SCA to search for
493 distinctive and common information in linked data. *Behavior Research Methods* **46**, 576–587 (2013)

494 18. Lofsted, T., Hoffman, D., Trygg, J.: Global, local and unique decomposition in OnPLS for multiblock data
495 analysis. *Analytica Chimica Acta* **791**, 13–24 (2012)

496 19. Måge, I., Smilde, A.K., van der Kloet, F.M.: Performance of methods that separate common and distinct
497 variation in multiple data blocks. *Journal of Chemometrics* **33**, 3085 (2019)

498 20. McCabe, S.D., Lin, D.Y., Love, M.I.: Consistency and overfitting of multi-omics methods on experimental data.
499 *Briefings in Bioinformatics* **21**, 1277–1284 (2020)

500 21. Lund, E., Dumeaux, V., Braaten, T., Hjartåker, A., Engeset, D., Skeie, G., Kumle, M.: Cohort profile: the
501 Norwegian Women and Cancer study: NOWAC - kvinner og kreft. *International Journal of Epidemiology* **37**,
502 36–4 (2008)

503 22. Zhu, M., Ghodsi, A.: Automatic dimensionality selection from the scree plot via the use of profile likelihood.
504 *Computational Statistics and Data Analysis* **51**, 918–930 (2006)

505 23. Jiang, M.: AJIVE Project. https://github.com/MeileiJiang/AJIVE_Project (2018).
506 https://github.com/MeileiJiang/AJIVE_Project

507 24. Carmichael, I.: ajive: Angle based Joint and Individual Variation Explained. https://github.com/idc9/r_jive
508 (2019). https://github.com/MeileiJiang/AJIVE_Project

509 25. Tang, T.M.: Integrated Principal Components Analysis (iPCA). <https://github.com/DataSlingsers/iPCA>
510 (2018). https://github.com/MeileiJiang/AJIVE_Project

511 26. Yanaihara, N., Caplen, N., Bowman, E., Seike, M., Kumamoto, K., Yi, M., Stephens, R.M., Okamoto, A.,
512 Yokota, J., Tanaka, T., Calin, G.A., Liu, C.-G., Croce, C.M., Harris, C.C.: Unique microRNA molecular profiles in
513 lung cancer diagnosis and prognosis. *Cancer Cell* **9**(3), 189–198 (2006). doi:[10.1016/j.ccr.2006.01.025](https://doi.org/10.1016/j.ccr.2006.01.025)

514 27. Hu, Y., Chen, G.: Pathogenic mechanisms of lung adenocarcinoma in smokers and non-smokers determined by
515 gene expression interrogation. *Oncology Letters* **10**, 1350–1370 (2015)

516 28. Zhang, Y., Breitling, L.P., Balavarca, Y., Hollecsek, B., Schöttker, B., Brenner, H.: Comparison and
517 combination of blood DNA methylation at smoking-associated genes and at lung cancer-related genes in

- 518 prediction of lung cancer mortality. *International Journal of Cancer* **139**(11), 2482–2492 (2016).
519 doi:10.1002/ijc.30374. <https://onlinelibrary.wiley.com/doi/pdf/10.1002/ijc.30374>
- 520 29. Baglietto, L., Ponzi, E., Haycock, P., Hodge, A., Assumma, M.B., Jung, C.H., Chung, J., Fasanelli, F., Guida,
521 F., Campanella, G., Chadeau-Hyam, M., Grankvist, K., Johansson, M., Ala, U., Provero, P., Wong, E.M., Joo,
522 J., English, D.R., Kazmi, N., Lund, E., Faltus, C., Kaaks, R., Risch, A., Barrdahl, M., Sandanger, T.M.,
523 Southey, M.C., Giles, G.G., Johansson, M., Vineis, P., Polidoro, S., Relton, C.L., Severi, G.: DNA methylation
524 changes measured in pre-diagnostic peripheral blood samples are associated with smoking and lung cancer risk.
525 *International Journal of Cancer* **140**, 50–61 (2017)
- 526 30. Heller, G., Weinzierl, M., Noll, C., Babinsky, V., Ziegler, B., Altenberger, C., Minichsdorfer, C., Lang, G.,
527 Döme, B., End-Pfützenreuter, A., Arns, B.-M., Grin, Y., Klepetko, W., Zielinski, C.C., Zöchbauer-Müller, S.:
528 Genome-wide mirna expression profiling identifies mir-9-3 and mir-193a as targets for dna methylation in
529 non-small cell lung cancers. *Clinical Cancer Research* **18**(6), 1619–1629 (2012).
530 doi:10.1158/1078-0432.CCR-11-2450. <https://clincancerres.aacrjournals.org/content/18/6/1619.full.pdf>
- 531 31. Sandanger, T.M., Haugdahl Nøst, T., Guida, F., Rylander, C., Campanella, G., Muller, D.C., van Dongen, J.,
532 Boomsma, D.I., Johansson, M., Vineis, P., Vermeulen, R., Lund, E., Chadeau-Hyam, M.: DNA methylation and
533 associated gene expression in blood prior to lung cancer diagnosis in the norwegian women and cancer cohort.
534 *Scientific Reports* **8**, 16714 (2018)
- 535 32. Zhang, Z., Liu, J., Kaur, M., Krantz, I.D.: Characterization of DNA methylation and its association with other
536 biological systems in lymphoblastoid cell lines. *Genomics* **99**(4), 209–219 (2012).
537 doi:10.1016/j.ygeno.2012.01.002
- 538 33. Ma, B., Wilker, E.H., Willis-Owen, S.A.G., Byun, H.-M., Wong, K.C.C., Motta, V., Baccarelli, A.A., Schwartz,
539 J., Cookson, W.O.C.M., Khabbaz, K., Mittleman, M.A., Moffatt, M.F., Liang, L.: Predicting DNA methylation
540 level across human tissues. *Nucleic Acids Research* **42**(6), 3515–3528 (2014). doi:10.1093/nar/gkt1380
- 541 34. Du, P., Zhang, X., Huang, C.C., Jafari, N., Kibbe, W.A., Hou, L., Lin, S.M.: Comparison of Beta-value and
542 M-value methods for quantifying methylation levels by microarray analysis. *BMC Bioinformatics* **11**, 587 (2010)
- 543 35. Jiang, R., Tang, W., Wu, X., Wenhui, F.: A random forest approach to the detection of epistatic interactions in
544 case-control studies. *BMC Bioinformatics* **10**, 65 (2009)
- 545 36. Fasanelli, F., Baglietto, L., Ponzi, E., Guida, F., Campanella, G., Johansson, M., Grankvist, K., Johansson, M.,
546 Assumma, M.B., Naccarati, A., Chadeau-Hyam, M., Ala, U., F., Kaaks, R., Risch, A., De Stavola, B., Hodge,
547 A., Giles, G.G., Southey, M.C., Relton, C.L., Haycock, P.C., Lund, E., Polidoro, S., Sandanger, T.M., Severi,
548 G., Vineis, P.: Hypomethylation of smoking-related genes is associated with future lung cancer in four
549 prospective cohorts. *Nature Communications* **6**, 10192 (2015)
- 550 37. Gayananova, I., Li, G.: Structural learning and integrative decomposition of multi-view data. *Biometrics* **75**,
551 1121–1132 (2019)

552 Figures

Figure 1: **aJIVE decomposition of the three omics sources**

The joint component dominates the individual and residual components for all three datasets.

Figure 2: **Joint and individual proportions of variance explained**

The joint component is prevalent on the individual and residual components for all three datasets.

Figure 3: **ROC curves from 10 fold cross validation**

a) reports the ROC curves and their AUCs for the prediction models on case vs control, b) reports the ROC curves and their AUCs for the prediction models on metastasis status

Figure 4: **Variable importance plot from random forest on case vs control**

First ten variables ranked by variable importance (in terms of mean Gini index) in the full model for case vs control

Figure 5: **Prediction models ROC from iPCA and aJIVE** a) reports the ROC curves and their AUCs for the prediction models on case vs control, b) reports the ROC curves and their AUCs for the prediction models on metastasis status

Comparison	Joint Rank	Ind Rank Methylation	Ind Rank mRNA	Ind Rank miRNA
methyl-mRNA	5	38	6	-
methyl-miRNA	2	41	-	7
mRNA-miRNA	3	-	8	6

Table 1: Joint and individual ranks obtained by pairwise aJIVE.

Model	Accuracy	OOB Classification Error	mean AUC
Joint and Individual (with covs)	0.81	26.53%	0.77
Joint (with covs)	0.68	33.67%	0.74
Patient covariates	0.65	38.27%	0.61
Individual (with covs)	0.77	30.61%	0.71
Non-integrative analysis (with covs)	0.50	38.05%	0.65

Table 2: Random forest diagnostics

Figures

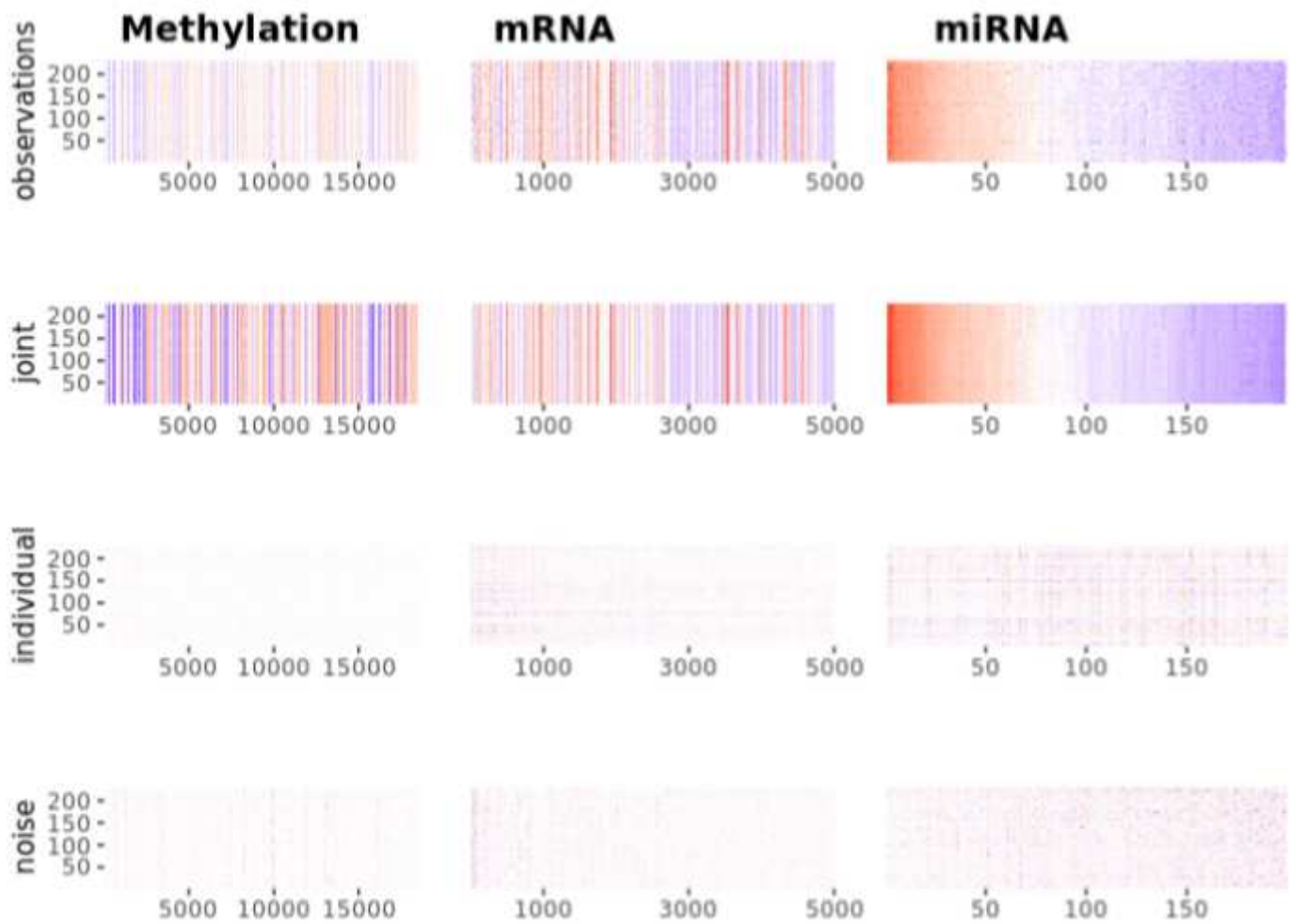


Figure 1

aJIVE decomposition of the three omics sources. The joint component dominates the individual and residual components for all three datasets.

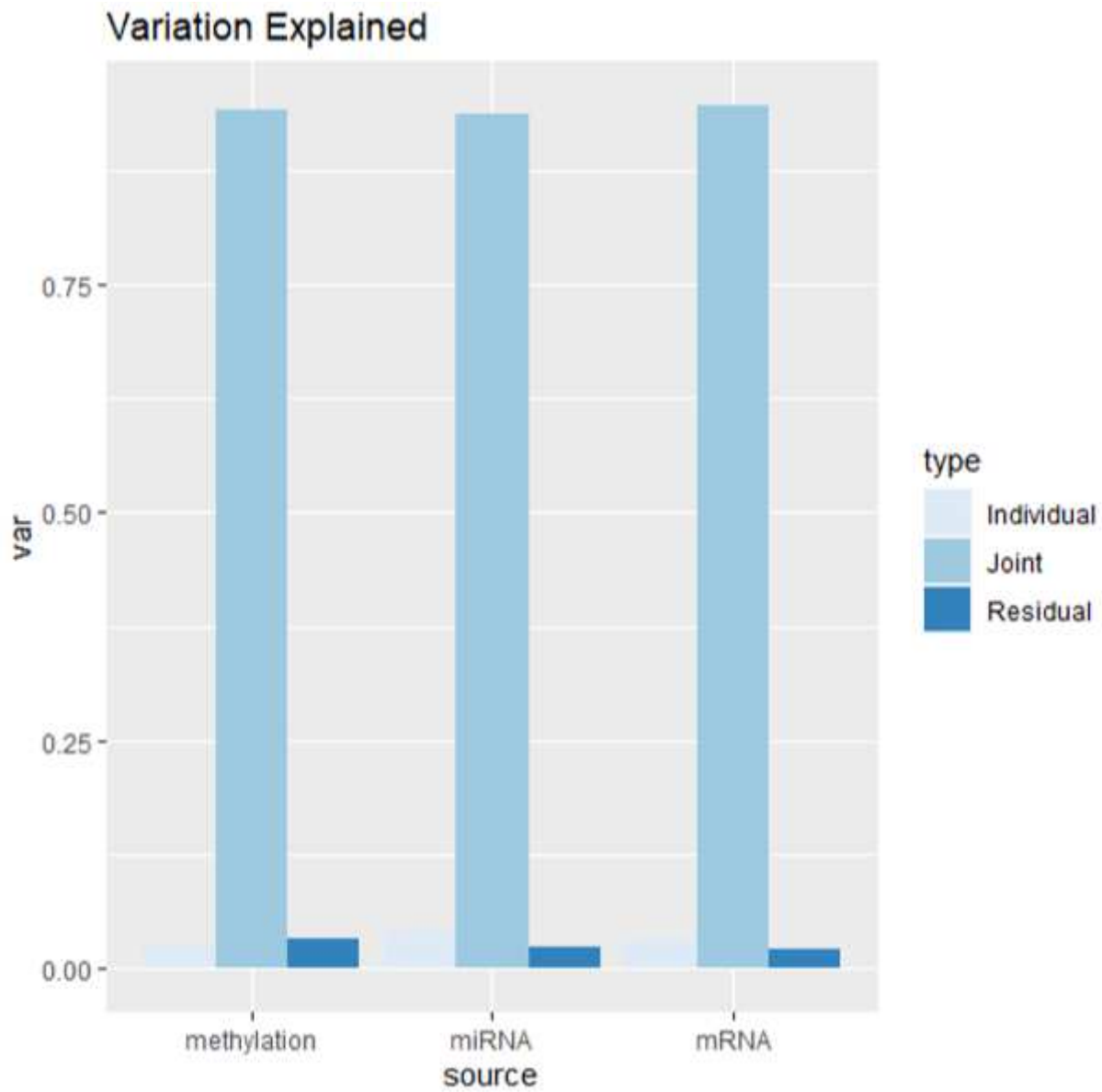


Figure 2

Joint and individual proportions of variance explained The joint component is prevalent on the individual and residual components for all three datasets.

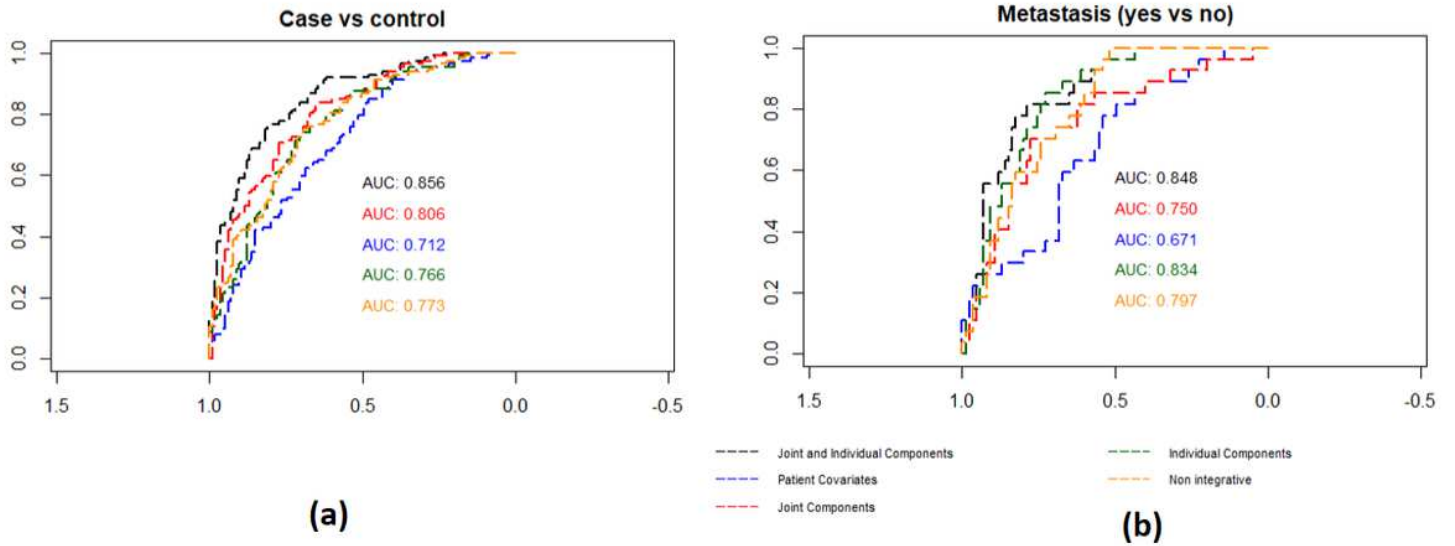


Figure 3

ROC curves from 10 fold cross validation a) reports the ROC curves and their AUCs for the prediction models on case vs control, b) reports the ROC curves and their AUCs for the prediction models on metastasis status

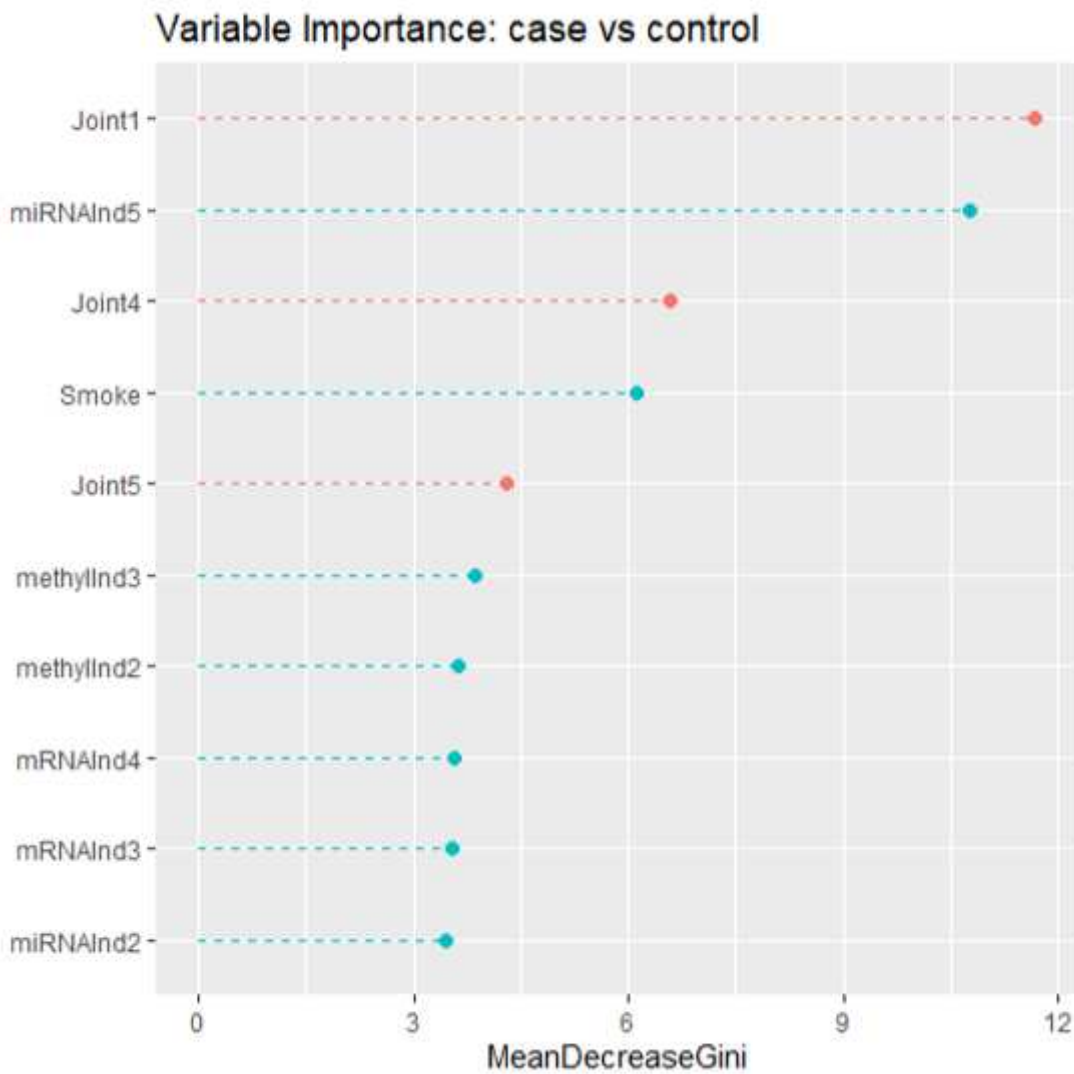


Figure 4

Variable importance plot from random forest on case vs control First ten variables ranked by variable importance (in terms of mean Gini index) in the full model for case vs control

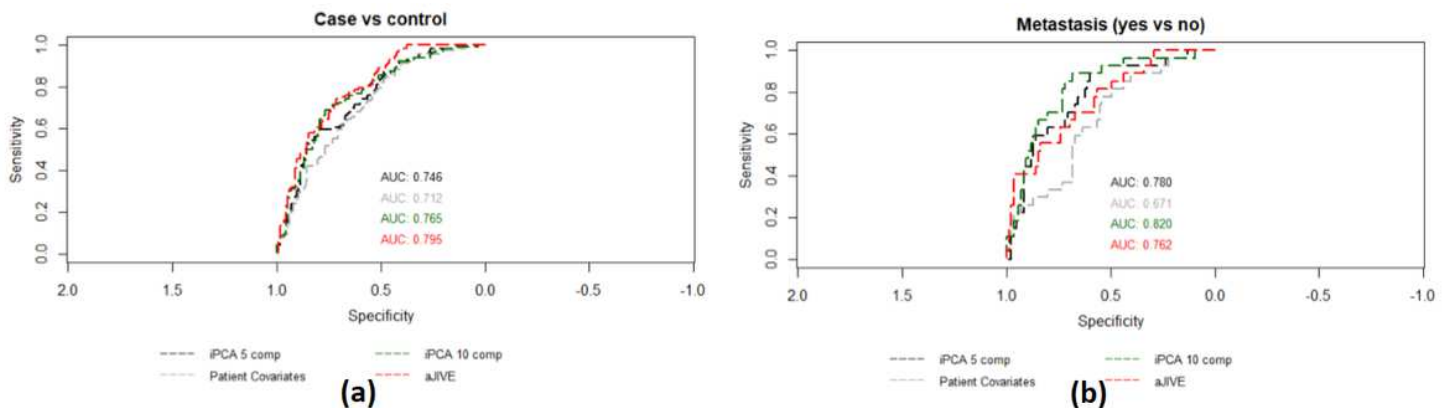


Figure 5

Prediction models ROC from iPCA and aJIVE a) reports the ROC curves and their AUCs for the prediction models on case vs control, b) reports the ROC curves and their AUCs for the prediction models on metastasis status