

# Integrating Multi-omics with Sparse Canonical Correlation Analysis

Michael Komodromos

`michael.komodromos19@ic.ac.uk`

## Abstract

We integrate radiomic and transcriptomic data from patients with ovarian cancer using sparse canonical correlation analysis (sCCA). We demonstrate integration yields prognostic models with greater predictive accuracy in comparison to current models. Further we examine network structures providing plausible relational pathways between genes and radiomic features.

**Code:** our study can be reproduced by running the code available at <https://github.com/mkomod/ovc>. Additionally, we provide an R package to perform sCCA found at <https://github.com/mkomod/rcca>.

## 1 Introduction

Multi-omics is a form of multi-view data, where several distinct feature sets are measured for the same samples; providing complementary information when characterising a biological object (Li et al., 2018). There are several types of “omics” data including: genomics, transcriptomics, proteomics, metabolomics etc.. A recent development in multi-omics is radiomics, which provide a set of quantitative features extracted from an image. Radiomics features include various first and higher order statistics, fractal and shape features characterising disease features not appreciated by the naked eye (Lu et al., 2019). Several studies have explored the integration of multi-omics data (Pittman et al., 2004; Hasin et al., 2017; Chaddad et al., 2019) and have demonstrated the increased predictive capacity of integration.

Canonical Correlation Analysis (CCA) is a method used to study relations between sets of covariates (Hotelling, 1936). CCA relies on the correlation structure between two groups of data to find linear combinations of covariates that maximise the correlation. CCA and its extensions have been successfully used in several disciplines, including: economics, neuroscience and computational biology, to name a few (Zhuang et al., 2020; Yamanishi et al., 2003). Furthermore, there have been several studies using canonical correlation based methods to integrate multi-omics (Witten and Tibshirani, 2009; Hong et al., 2013; Shi et al., 2019; Rodosthenous et al., 2020). For instance, Shi et al. demonstrated that CCA based methods can be used to identify relevant miRNA-mRNA pathways.

In recent years there have been several extensions of CCA (Uurtio et al., 2017; Li et al., 2018). For instance, kernel CCA and deep CCA allow for non-linear relationships between groups of covariates to be discovered through the use (or construction) of a feature map (Andrew et al., 2013). Further, regularised forms of CCA known as sparse CCA (sCCA) have been proposed to tackle high-dimensional datasets, where there are often more covariates than observations (Witten and Tibshirani, 2009; Suo, 2018; Rodosthenous et al., 2020). Finally, CCA has also been framed in a probabilistic perspective allowing for Bayesian interpretations and correspondingly Bayesian equivalents to the mentioned extensions (Bach and Jordan, 2005).

Within our study we consider the integration of radiomic and transcriptomic data through sparse CCA; examining both predictive and exploratory aspects. Regarding the predictive performance we compare our results to current prognostic models. Further, we explore the shared structure between these datasets using network based methods.

## 2 Methods

### 2.1 Canonical Correlation Analysis (CCA)

CCA is concerned with finding pairs of linear projections of two views that are maximally correlated. Let  $(\mathbf{X}_1, \mathbf{X}_2) \in \mathbb{R}^{p_1} \times \mathbb{R}^{p_2}$  be random vectors with covariances  $(\Sigma_{11}, \Sigma_{22})$  and cross-covariance  $\Sigma_{12} = \Sigma_{21}^\top$ . Consider a pair of linear projections  $(w_1^\top \mathbf{X}_1, w_2^\top \mathbf{X}_2)$ , we are interested in finding the pair of vectors  $(w_1^*, w_2^*)$  that

maximise the correlation between the linear projections  $w_1^\top \mathbf{X}_1, w_2^\top \mathbf{X}_2$ . Formally,

$$(w_1^*, w_2^*) = \underset{w_1, w_2}{\operatorname{argmax}} \operatorname{corr}(w_1^\top \mathbf{X}_1, w_2^\top \mathbf{X}_2)$$

which is equivalent to solving

$$\underset{w_1, w_2}{\operatorname{maximise}} w_1^\top \Sigma_{12} w_2, \quad \text{subject to } w_1^\top \Sigma_{11} w_1 = w_2^\top \Sigma_{22} w_2 = 1. \quad (1)$$

Notably, we can find up to  $\min(p_1, p_2)$  pairs of vectors denoted as  $(w_1^{(i)}, w_2^{(i)})$  for  $i = 1, \dots, \min(p_1, p_2)$ . Let  $A_1 \in \mathbb{R}^{p_1 \times m}$  and  $A_2 \in \mathbb{R}^{p_2 \times m}$  where the columns of  $A_1$  and  $A_2$  correspond to the first  $m$  vectors  $w_1^{(i)}$  and  $w_2^{(i)}$  for  $m \leq \min(p_1, p_2)$ . Finding the first  $m$  vector pairs  $(w_1^{(i)}, w_2^{(i)})$  corresponds to solving

$$\underset{A_1, A_2}{\operatorname{maximise}} \operatorname{tr}(A_1^\top \Sigma_{12} A_2) \quad \text{subject to } A_1^\top \Sigma_{11} A_1 = A_2^\top \Sigma_{22} A_2 = I. \quad (2)$$

where  $\operatorname{tr}(\cdot)$  is the trace. Let  $T := \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2}$ , if we suppose  $\Sigma_{11}$  and  $\Sigma_{22}$  are non-singular, then the singular value decomposition of  $T$  can be written in the form

$$T = (u_1, \dots, u_m) D (v_1, \dots, v_k)^\top \quad (3)$$

where  $u_i \in \mathbb{R}^{p_1}$ ,  $v_i \in \mathbb{R}^{p_2}$  and  $D = \operatorname{diag}(\lambda_1, \dots, \lambda_m)$ , then the  $i$ th *canonical correlation vectors* are given by

$$w_1^{(i)} = \Sigma_{11}^{-1/2} u_i \quad (4)$$

$$w_2^{(i)} = \Sigma_{22}^{-1/2} v_i, \quad (5)$$

and *canonical correlation coefficients*  $\rho_i = \lambda_i$  for  $i = 1, \dots, m$  (Mardia et al., 1979).

## 2.2 Sparse Canonical Correlation Analysis (sCCA)

Traditional CCA breaks down in a high dimensional setting. For example, consider (centered) observations  $Y \in \mathbb{R}^{n \times p}$  where  $p > n$ , the covariance matrix is given by  $\Sigma = Y^\top Y \in \mathbb{R}^{p \times p}$ . We notice,  $\operatorname{rank}(\Sigma) \leq n - 1 < p$ , hence  $\Sigma$  is singular and we are unable to evaluate  $T$ . sCCA adapts the objective function presented in Equation (1) so we can perform correlation analysis in a high dimensional setting.

First, we present a generalised framework for regularised CCA, and then present methods by (Witten and Tibshirani, 2009) and (Suo, 2018). Contrary to our description of CCA, we consider matrices of observations, denoted as  $X$ . Let  $X_1 \in \mathbb{R}^{n \times p_1}$

and  $X_2 \in \mathbb{R}^{n \times p_2}$  be matrices of observations where  $p \gg n$ . Let  $w_1 \in \mathbb{R}^{p_1}$  and  $w_2 \in \mathbb{R}^{p_2}$  and let  $r_1 : w_1 \mapsto \mathbb{R}$  and  $r_2 : w_2 \mapsto \mathbb{R}$ . Then, the generalised form for the regularised CCA problem is given by

$$\underset{w_1, w_2}{\text{maximise}} \quad \text{corr}(X_1 w_1, X_2 w_2) - r_1(w_1) - r_2(w_2). \quad (6)$$

An efficient algorithm to find sparse canonical vectors was proposed by Witten and Tibshirani (2009), under which they set  $r_1$  and  $r_2$  to be the  $\ell_1$ -norm and solve

$$\begin{aligned} &\underset{w_1, w_2}{\text{maximise}} \quad \text{corr}(X_1 w_1, X_2 w_2) - \lambda_1 \|w_1\|_1 - \lambda_2 \|w_2\|_1 \\ &\text{subject to} \quad \|w_1\|_2 \leq 1, \|w_2\|_2 \leq 1, \end{aligned} \quad (7)$$

where  $\|\cdot\|_1$  denotes the  $\ell_1$ -norm,  $\|\cdot\|_2$  denotes the  $\ell_2$ -norm and  $\lambda_1, \lambda_2 \in \mathbb{R}_+$  are regularisation hyperparameters. Under the above formulation the canonical vectors can be inaccurate and non-sparse (Suo, 2018; Rodosthenous et al., 2020); hence Suo relax the constraints in Equation (7) to be

$$\begin{aligned} &\underset{w_1, w_2}{\text{maximise}} \quad \text{corr}(X_1 w_1, X_2 w_2) - \lambda_1 \|w_1\|_1 - \lambda_2 \|w_2\|_1. \\ &\text{subject to} \quad \text{var}(X_1 w_1) \leq 1, \text{var}(X_2 w_2) \leq 1. \end{aligned} \quad (8)$$

The resulting problem is biconvex, in other words, fixing  $w_1$  makes the problem convex with respect to  $w_2$  and fixing  $w_2$  makes the problem convex with respect to  $w_1$ . A description of the algorithm used to solve Equation (8) is presented in (Suo, 2018).

### 2.2.1 Multiple canonical vector pairs

As with conventional CCA we are able to find multiple pairs of canonical vectors. However, we do so in an iterative manner, by first solving Equation (8) and then solving again where instead of  $X_1$  and  $X_2$  we use

$$\bar{X}_1 = \begin{pmatrix} X_1 \\ W_1^\top X_1^\top X_1 \\ W_2^\top X_2^\top X_1 \end{pmatrix}, \quad \bar{X}_2 = \begin{pmatrix} X_2 \\ W_2^\top X_2^\top X_2 \\ W_1^\top X_1^\top X_2 \end{pmatrix} \quad (9)$$

where  $W_k = \left(w_k^{(i)}\right)_{i=1}^{r-1}$  for  $k = 1, 2$  are matrices with the  $(r-1)$  previous canonical vectors (Rodosthenous et al., 2020). Solving yields the  $r$ th pair of canonical vectors, noting we are able to obtain up to  $\min(p_1, p_2)$  pairs.

### 2.2.2 Tuning the regularisation hyperparameters

There are a number of methods for tuning the regularisation hyperparameters. As our dataset is small we opt to use a permutation based tuning method proposed by (Witten and Tibshirani, 2009). Witten and Tibshirani’s method is outlined in Algorithm 1. Noting, we have adapted the method to optimise Suo’s method for sCCA.

---

**Algorithm 1** Permutation validation for  $(\lambda_1, \lambda_2)$ .

---

```

1: for  $(\lambda_1, \lambda_2)_j \in \Lambda$  do
2:   Solve Equation (8) for  $(w_1^*, w_2^*)$  for  $X_1$  and  $X_2$ , taking  $(\lambda_1, \lambda_2) = (\lambda_1, \lambda_2)_j$ .
3:   Compute  $d_j = \text{corr}(X_1 w_1^*, X_2 w_2^*)$ .
4:   for  $i \in \{1, \dots, B\}$  do
5:     Permute the rows of  $X_1$ ; constructing  $X_1^i$ .
6:     Solve Equation (8)  $(w_1^i, w_2^i)$  for  $X_1^i$  and  $X_2$ , taking  $(\lambda_1, \lambda_2) = (\lambda_1, \lambda_2)_j$ .
7:     Compute  $d_j^i = \text{corr}(X_1^i w_1^i, X_2 w_2^i)$ ,
8:   end for
9:   Compute  $p_j = \frac{1}{B} \sum_{i=1}^B 1(d_j^i \geq d_j)$ .
10: end for
11: return  $(\lambda_1, \lambda_2)_j$  that gives the smallest  $p_j$ .
```

---

### 2.3 Semi-supervised sparse CCA

Semi-supervised sparse CCA (SS-CCA) is a method used to select features of  $X_1$  and  $X_2$  related to a response  $Y \in \mathbb{R}^n$ . The objective function of SS-CCA is given as

$$\underset{w_1, w_2}{\text{maximise}} \quad \text{corr}(X_1 w_1, X_2 w_2) - \lambda_1 \|w_1\|_1 - \lambda_2 \|w_2\|_1 \quad (10)$$

$$\text{subject to} \quad \text{var}(X_1 w_1) \leq 1, \quad \text{var}(X_2 w_2) \leq 1,$$

$$w_{1i} = 0 \quad \forall i \in Q_1, \quad w_{2j} = 0 \quad \forall j \in Q_2.$$

where  $Q_1$  and  $Q_2$  are the set of indices of features of  $X_1$  and  $X_2$  that are least correlated with the response  $Y$  (Witten and Tibshirani, 2009). Generalising the

criteria for  $Q_1$  and  $Q_2$ , we have

$$Q_1 = \{i : f(X_{1i}, Y) \leq \gamma_1, i = 1, \dots, p_1\}, \quad (11)$$

$$Q_2 = \{j : f(X_{2j}, Y) \leq \gamma_2, j = 1, \dots, p_2\} \quad (12)$$

where  $\gamma_1$  and  $\gamma_2$  are thresholds controlling the elements in  $Q_1$  and  $Q_2$ , and  $f$  is a function to evaluate the relatedness of a feature of  $X$  to  $Y$ . For example, when considering overall survival,  $f$  can return the Cox's statistic. Notably,  $\gamma_1$  and  $\gamma_2$  are hyperparameters that need to be tuned.

## 2.4 Proportional Hazards Model

The proportional hazards model (PHM) is used to model time to failure events. Let  $T$  denote a random variable representing time to failure, with density  $f(t)$  and survivor function  $S(t) = \mathbb{P}(T > t)$ . The hazard rate, the rate of failure at time  $t$ , is given by

$$\lambda(t) = \frac{f(t)}{S(t)}, \quad t > 0. \quad (13)$$

The PHM is based on the assumption that the hazard rate is a product of a baseline hazard rate,  $\lambda_0(t)$ , and a positive functional term  $\psi(x; \beta)$ , formally,

$$\lambda(t, x) = \lambda_0(t)\psi(x; \beta), \quad (14)$$

for covariates  $x \in \mathbb{R}^p$  and unknown parameters  $\beta \in \mathbb{R}^p$ . A number of forms have been proposed for  $\psi(x; \beta)$ , the most common being  $\psi(x; \beta) = \exp(\beta^\top x)$ . We note, an implication of the form of Equation (14) is that the ratio of any two samples at time  $t$  is constant, i.e.  $\lambda(t, x_1) \propto \lambda(t, x_2)$  (Cox, 1972). Parameter values  $\beta$  can be estimated using the partial log-likelihood,  $\ell(\beta)$  (Cox, 1975).

### 2.4.1 Evaluating proportional hazard models

The most widely used statistic to evaluate the performance of PHMs is concordance. Concordance is measure of the discriminatory power of a model. Let  $t_i$  and  $\hat{t}_i$  be observed and predicted failure times for observation  $i$ , then the concordance,

$$c = \mathbb{P}(\hat{t}_i > \hat{t}_j | t_i > t_j), \quad (15)$$

for a pair of observations  $i, j$ . A popular estimator for the concordance is Harrell's  $c$ -index, given as

$$\hat{c} = \frac{\sum \sum_{i < j} \mathbb{I}(t_i < t_j) \mathbb{I}(\hat{\beta}^\top x_i > \hat{\beta}^\top x_j) \delta_i + \mathbb{I}(t_j < t_i) \mathbb{I}(\hat{\beta}^\top x_j > \hat{\beta}^\top x_i) \delta_j}{\sum \sum_{i < j} \mathbb{I}(t_i < t_j) \delta_i + \mathbb{I}(t_j < t_i) \delta_j}, \quad (16)$$

where  $\mathbb{I}(\cdot)$  is the indicator function and  $\delta_i = 1$  if the observation is uncensored and 0 otherwise. Notably, the  $c$ -index is sensitive to the degree of censoring and tends to be higher when the degree of censoring is higher. As a result, an alternate estimator proposed by Gonen and Heller is used for comparison, given as

$$\hat{k} = \frac{1}{n(n-1)} \sum \sum_{i < j} \left( \frac{\mathbb{I}(\hat{\beta}^\top x_{ji} < 0)}{1 + \exp(\hat{\beta}^\top x_{ji})} + \frac{\mathbb{I}(\hat{\beta}^\top x_{ij} < 0)}{1 + \exp(\hat{\beta}^\top x_{ij})} \right), \quad (17)$$

where  $x_{ij}$  represents pairwise difference  $x_i - x_j$ . As  $\hat{k}$  is determined through the partial likelihood estimate,  $\hat{\beta}$ , and the effect of censoring on the bias of  $\hat{\beta}$  is negligible; we obtain a robust estimator for the concordance probability (Gonen and Heller, 2005). Notably, we are able to obtain an estimate for the standard error of  $\hat{k}$  via a smooth approximation for  $\hat{k}$ , details are presented in (Gonen and Heller, 2005).

## 2.5 Networks

Networks can be used to understand the relationships within groups of covariates. To construct networks we use a method proposed by Shi et al. referred to as SmCCNet. SmCCNet uses a similarity matrix constructed by robust canonical vector pairs, details are outlined in Algorithm 2. We note, unlike the original algorithm we use Suo's method to performs sCCA. Following the construction of  $S$ , hierarchical clustering is performed and an appropriate level is set to cut the dendrogram. Cliques containing features from both  $X_1$  and  $X_2$  are returned and used to construct networks. Full details are presented in (Shi et al., 2019).

---

**Algorithm 2** SmCCNet similarity matrix.

---

```
1: for  $i \in \{1, \dots, B\}$  do
2:   Subsample columns of  $X_1$  and  $X_2$ , we denote these matrices as  $X'_1$  and  $X'_2$ .
3:   Compute  $(w'_1, w'_2)$  by solving Equation (8) for  $X'_1$  and  $X'_2$ 
4:   Compute  $A_i = w'_1 \otimes w'_2$ .
5: end for
6: Compute  $A$  the element-wise average over the entires in  $A_1, \dots, A_B$ 
7: Normalise  $A$  by  $\max(A)$ .
8: return  $S = 1 - A$ 
```

---

### 3 Results

Data for this study consists of clinical, radiomic and genomic data available from the cancer genome atlas (TCGA). Radiomics data is comprised of features relating to: size, intensity, texture and wavelet decompositions of contrast enhanced CT scans (Lu et al., 2019). The radiomics data for this study were produced by Lu et al., the remaining datasets are summarised in Table 1. Notably, both the radiomics and mRNA datasets have more features than samples.

Data type	Platform	Samples	Features
Clinical		630	19
Radiomics	TexLab 2.0	71	658
mRNA expression	Affymetrix U133	593	12,043

Table 1: Summary of data used.

We merged the datasets to construct our multi-view dataset, upon merging there were  $n = 68$  samples, 50 of which had right censored overall survival times (days). Further, prior to our analysis we pre-selected covariates from the Radiomic and mRNA expression datasets by fitting univariate proportional hazards models and selecting covariates with  $p$ -values less than  $\alpha = 0.05$ . Variable selection left  $p_1 = 524$  features in our mRNA expression dataset and  $p_2 = 32$  features in our radiomics dataset. Notably, applying sCCA to our filtered datasets is equivalent to conducting



semi-supervised sparse CCA.

### 3.1 Canonical Correlation Analysis

We applied Suo’s method for sparse CCA to our Radiomics and mRNA datasets, centering and standardising beforehand. We tuned the regularisation hyperparameters,  $\lambda_1$  and  $\lambda_2$ , in Equation (10) using permutation validation (Algorithm 1) taking  $\Lambda = \{0.4, 0.5, \dots, 2\} \times \{0.1, 0.2, \dots, 3\}$  and  $B = 1000$ . The loss landscape is presented in Figure 1, notably a minimum  $p$ -value of 0.018 occurs at (1.1, 1.6).

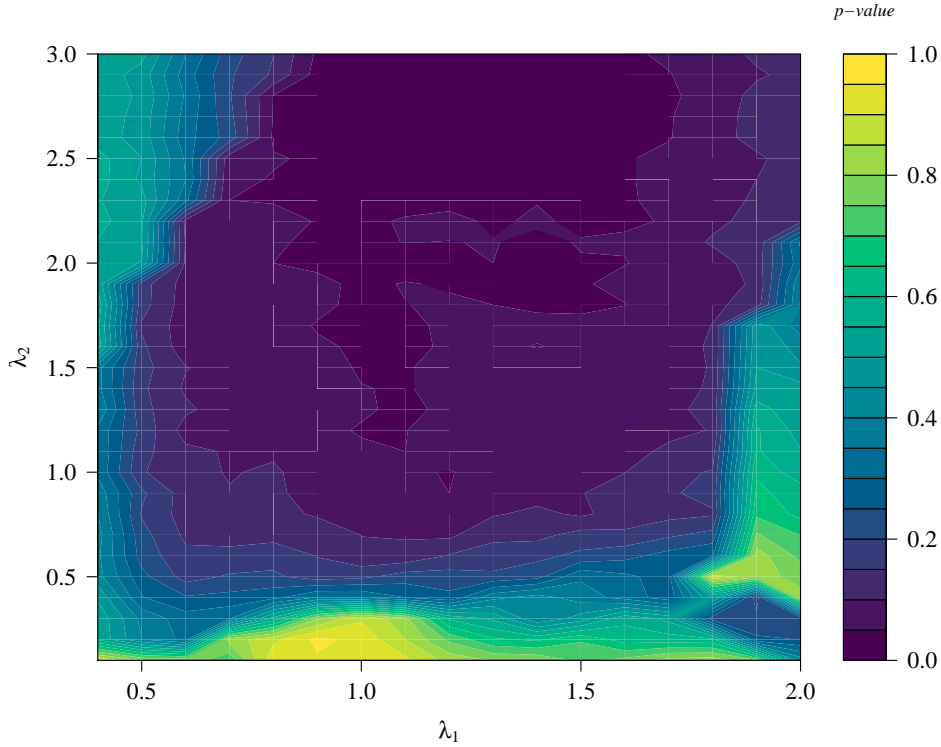


Figure 1: Loss landscape for a grid of values for  $\lambda_1$  and  $\lambda_2$ . We note, as  $\lambda_1 \rightarrow 0$  and  $\lambda_2 \rightarrow 0$  we impose less penalisation on  $w_1$  and  $w_2$ , leading to non-sparse results. As such, we avoid choosing regularisation parameters close to the origin.

Taking  $\lambda_1 = 1.1$  and  $\lambda_2 = 1.6$ , we obtain the first three canonical vector pairs by solving Equation (10). We present the canonical vector pairs in Figure 2. Notably, the vector pairs for the mRNA expression data ( $w_1^{(i)}$ ) are unstable, arising from the small sample size of the dataset (Suo, 2018). Further, we notice there are only two non-zero features in the first canonical vector of the radiomics dataset ( $w_2^{(1)}$ ), and

three for the remaining radiomics canonical vectors.

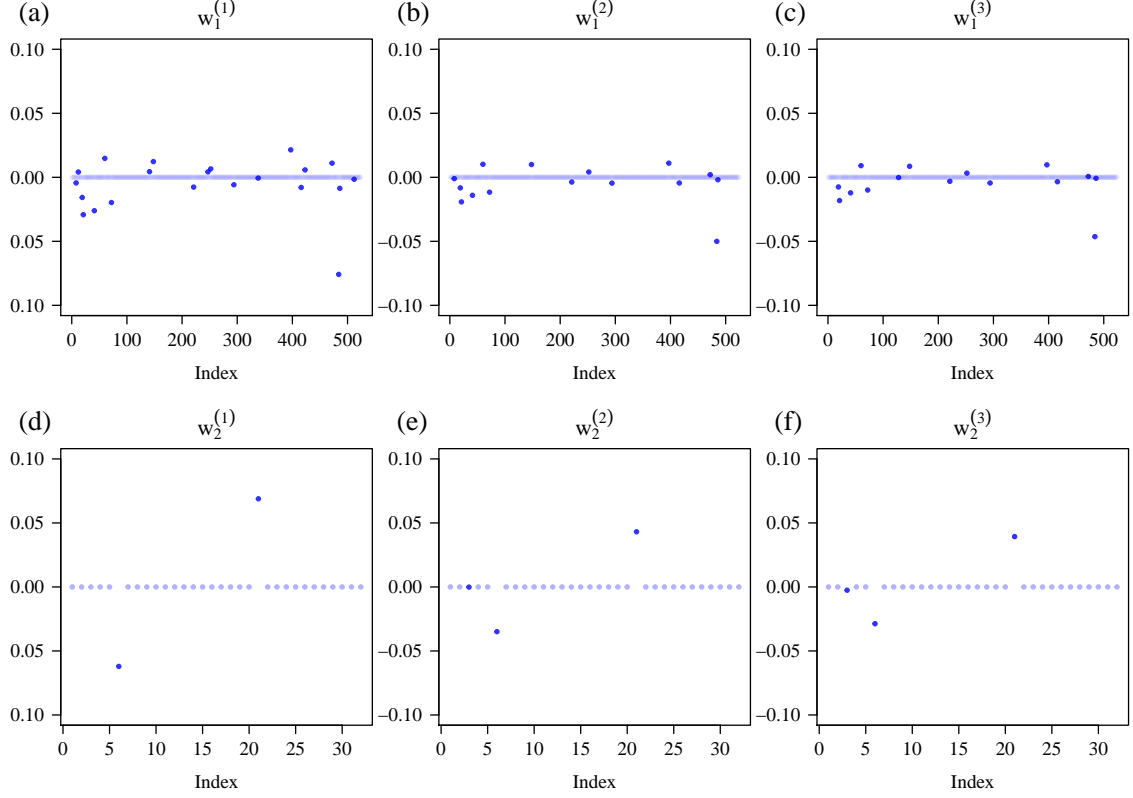


Figure 2: First three sparse canonical vectors for the mRNA expression dataset (a)-(c) and Radiomics dataset (d)-(f). For clarity elements of  $w^{(i)}$  with a value of 0 have a lighter shade. We notice a combination of few covariates are needed to maximise the correlation between the mRNA dataset and Radiomics dataset.

### 3.2 Models

To evaluate the canonical vector pairs we fit multivariate proportional hazards models taking overall survival as our response. We took age (dichotomised at 60 years), stage and projections of radiomics or mRNA data, where projections are given by  $Xw^{(i)}$ , as our predictors. As a baseline to compare our models against we consider a model adapted from Lu et al., taking age, stage and the radiomic prognostic vector (RPV) as the predictors, denoted as (M0). Our baseline model gave  $\hat{c} = 0.672$  (0.0978) and  $\hat{k} = 0.707$  (0.0510) where standard errors are given in parentheses.

Notably approximately 74% (50/68) samples in our dataset were censored. To

mitigate the issue of censoring in evaluating our models performance we compute Gonen and Heller’s estimate for the concordance index,  $\hat{k}$ . The performance of our models is presented in Table 2. Notably, all models constructed using mRNA projections had  $\hat{c}$  and  $\hat{k}$  greater than our baseline model. Further, all models had  $\hat{k}$  greater than our baseline model.

Projection	Radiomics		mRNA	
	$\hat{c}$	$\hat{k}$	$\hat{c}$	$\hat{k}$
1	0.631 (0.100)	0.709 (0.234)	0.713 (0.099)	0.716 (0.223)
2	0.631 (0.100)	0.710 (0.234)	0.733 (0.099)	0.715 (0.218)
3	0.631 (0.100)	0.710 (0.234)	0.735 (0.099)	0.715 (0.232)

Table 2: Concordance (standard error) of models. The projections included in the models are listed in the leftmost column.

Further, we constructed a model (M1) using age, stage, the first mRNA projection and the second radiomics project as predictors. The combination of these features yielded a model with  $\hat{c} = 0.764$  (0.072) and  $\hat{k} = 0.761$  (0.048).

### 3.3 Networks

To examine relationships between radiomic and mRNA expression features we constructed networks using SmCCNet. We present our results in Figure 3.

## 4 Discussion

Our results highlight the predictive and exploratory utility of sCCA, in particular we have shown the subspaces found by sCCA are as good or better than our baseline comparison. Furthermore, our method of discovering linear combinations of covariates is semi-supervised, i.e. covariates are preselected based on some criteria and canonical covariates are found without knowledge of the response. In comparison, in our baseline model the RPV is constructed through a supervised approach, where the RPV is a linear combination of features found using an  $\ell_1$  regularised pro-

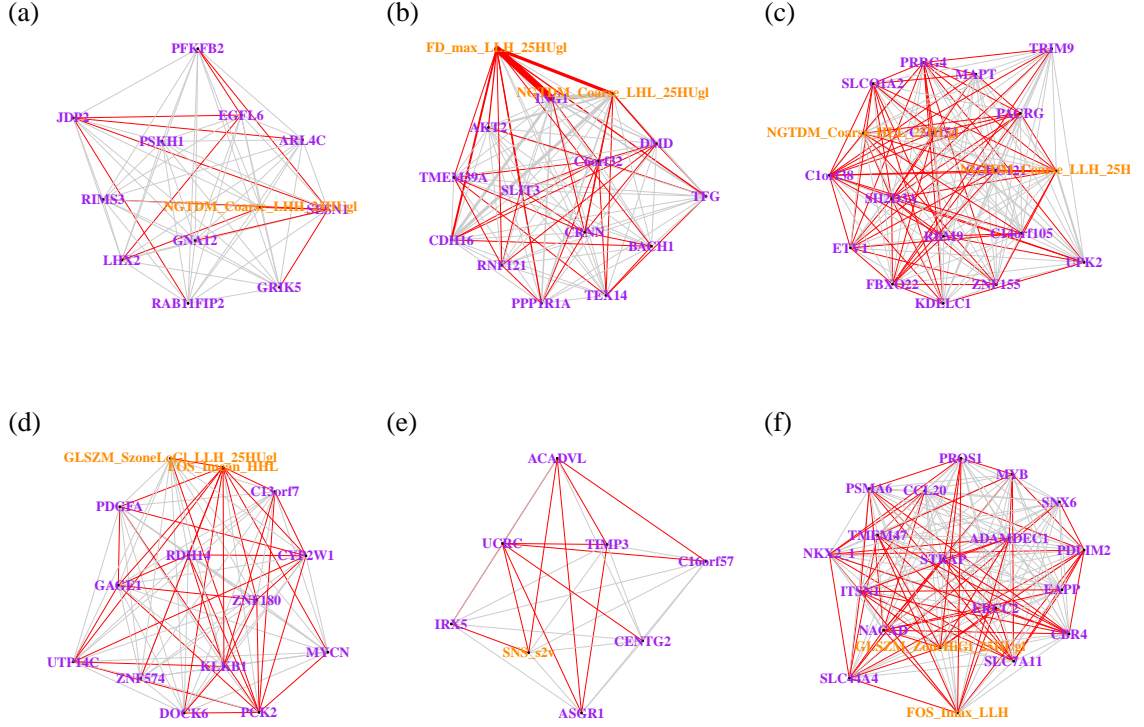


Figure 3: Networks constructed using SmCCNet. Radiomics features are coloured orange and genes are coloured purple. Connective lines indicate the correlation between the features with red lines indicating negative correlation and grey positive.

portional hazards model. Hence, not only does sCCA perform better than current models, subspaces are found without knowledge of the response. However, if prediction is our ultimate goal then a PHM comprised of age, stage and the projection of mRNA data onto the first principle component yields a model with  $\hat{c} = 0.86$  (0.046) and  $\hat{k} = 0.795$  (0.039). Outperforming all models in Table 2, M0 and M1. The difference in model performance suggests either the shared structure between the radiomics and the mRNA data is insubstantial (from a predictive perspective), or the signal to noise ratio is inhibiting the discovery of low dimensional structure.

Further, the interpretation of SmCCNet results requires caution. The networks formed are based on the outer product of the absolute canonical vector pairs; hence groups constructed are based on features with large absolute values. Furthermore, network pathways (connective lines) are constructed from the correlation matrix

between the mRNA and radiomics data. As estimation of the correlation matrix suffers in a high-dimensional setting, erroneous conclusions may be drawn from these pathways. As a results, care is needed when interpreting the results and networks should be used as a primer to direct further research rather than provide conclusive relations between mRNA and radiomics features.

From a methodological standpoint care has been taken to limit external input. However as shown by Figure 2, the permutation validation scheme for  $\lambda_1$  and  $\lambda_2$  yields results where the mRNA canonical vectors are unstable. We can correct for instability by increasing the value of  $\lambda_1$ . Suggesting, different validation schemes may need to be explored and their output examined; a comparison of these schemes may prove to be an insightful avenue for future research.

## 5 Conclusions

We have demonstrated by integrating radiomics and mRNA expression data we can increase the predictive accuracy of prognostic models. In particular, we have shown a combination of projections obtained from semi-supervised sparse canonical correlation analysis results in a model that outperforms current prognostic models derived from radiomics features. Further, we have uncovered pathways between radiomic and transcriptomics features that may prove insightful under further examination.

## References

- G. Andrew, R. Arora, J. Bilmes, and K. Livescu. Deep canonical correlation analysis. *30th Int. Conf. Mach. Learn. ICML 2013*, 28(3):2284–2292, 2013.
- F. R. Bach and M. I. Jordan. A Probabilistic Interpretation of Canonical Correlation Analysis. *Dept. Stat. Univ. California, Berkeley, CA, Tech. Rep*, pages 1–11, 2005.
- A. Chaddad, P. Daniel, S. Sabri, C. Desrosiers, and B. Abdulkarim. Integration of radiomic and multi-omic analyses predicts survival of newly diagnosed IDH1 wild-type glioblastoma. *Cancers (Basel)*., 11(8):1–16, 2019.
- D. R. Cox. Regression Models and Life-Tables. *J. R. Stat. Soc. Ser. B*, 34(2):187–220, Feb 1972.
- D. R. Cox. Partial Likelihood. *Biometrika*, 62(2):269–276, 1975.
- M. Gonen and G. Heller. Concordance probability and discriminatory power in proportional hazards regression. *Biometrika*, 92(4):965–970, 2005.
- Y. Hasin, M. Seldin, and A. Lusis. Multi-omics approaches to disease. *Genome Biol.*, 18(1):1–15, 2017.
- S. Hong, X. Chen, L. Jin, and M. Xiong. Canonical correlation analysis for RNA-seq co-expression networks. *Nucleic Acids Res.*, 41(8):1–15, 2013.
- H. Hotelling. Relations Between Two Sets of Variates. *Biometrika*, 28(3/4):321, 1936.
- Y. Li, F. X. Wu, and A. Ngom. A review on machine learning principles for multi-view biological data integration. *Brief. Bioinform.*, 19(2):325–340, 2018.
- H. Lu, M. Arshad, A. Thornton, G. Avesani, P. Cunnea, E. Curry, F. Kanavati, J. Liang, K. Nixon, S. T. Williams, M. A. Hassan, D. D. Bowtell, H. Gabra, C. Fotopoulou, A. Rockall, and E. O. Aboagye. A mathematical-descriptor of tumor-mesoscopic-structure from computed-tomography images annotates prognostic-

- and molecular-phenotypes of epithelial ovarian cancer. *Nat. Commun.*, 10(1): 1–11, 2019.
- K. V. Mardia, J. Kent, and J. Bibby. *Multivariate analysis*. Academic Press, 1979. ISBN 9780124712522.
- J. Pittman, E. Huang, H. Dressman, C. F. Horng, S. H. Cheng, M. H. Tsou, C. M. Chen, A. Bild, E. S. Iversen, A. T. Huang, J. R. Nevins, and M. West. Integrated modeling of clinical and gene expression information for personalized prediction of disease outcomes. *Proc. Natl. Acad. Sci. U. S. A.*, 101(22):8431–8436, 2004.
- T. Rodosthenous, V. Shahrezaei, and M. Evangelou. Integrating multi-OMICS data through sparse canonical correlation analysis for the prediction of complex traits: a comparison study. *Bioinformatics*, 36(17):4616–4625, 2020.
- W. J. Shi, Y. Zhuang, P. H. Russell, B. D. Hobbs, M. M. Parker, P. J. Castaldi, P. Rudra, B. Vestal, C. P. Hersh, L. M. Saba, and K. Kechris. Unsupervised discovery of phenotype-specific multi-omics networks. *Bioinformatics*, 35(21):4336–4343, 2019.
- X. Suo. *Topics In High-Dimensional Statistical Learning*. PhD thesis, Stanford, 2018.
- V. Uurtio, J. M. Monteiro, J. Kandola, J. Shawe-Taylor, D. Fernandez-Reyes, and J. Rousu. A tutorial on canonical correlation methods. *arXiv*, 50(6), 2017.
- D. M. Witten and R. J. Tibshirani. Extensions of sparse canonical correlation analysis with applications to genomic data. *Stat. Appl. Genet. Mol. Biol.*, 8(1): 1–27, 2009.
- Y. Yamanishi, J. P. Vert, A. Nakaya, and M. Kanehisa. Extraction of correlated gene clusters from multiple genomic data by generalized kernel canonical correlation analysis. *Bioinformatics*, 19 Suppl 1:323–330, 2003.
- X. Zhuang, Z. Yang, and D. Cordes. A technical review of canonical correlation analysis for neuroscience applications. *Hum. Brain Mapp.*, 41(13):3807–3833, 2020.