

Regression Analysis of Multi-view Microbiome Data

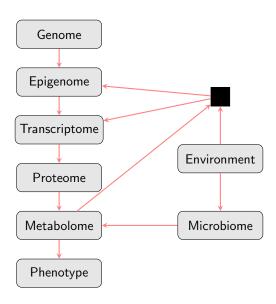
Jing Ma

Division of Public Health Sciences Fred Hutchinson Cancer Research Center

> 21 February 2022 IMSI Workshop

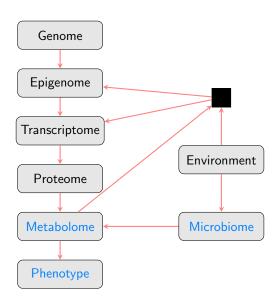
Systems Biology





Systems Biology





Outline



Motivation

Regression with One Dataset

Regression with Two Datasets

Future Work

BactoCARB Study



A randomized crossover study of 80 subjects aimed at the effects of high/low glycemic load on a variety of biomarkers.

- $X^{(1)} = 144$ metabolites & bile acids (blood)
- $X^{(2)} = 134$ microbial genus abundances (stool)
- y = enterolactone 'ENL' (urine), a product jointly produced by microbes and metabolites

Problem of Interest: which metabolites and bacteria are $\underline{jointly}$ associated with y?

BactoCARB Study



Maybe introduce a slide on how the data look like

Microbial Metabolism



etc ...

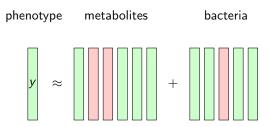


https://lpi.oregonstate.edu/mic/dietary-factors/phytochemicals/lignans

more is better!

Today's Goal





Find correlated metabolites and bacteria jointly associated with y

Outline



Motivation

Regression with One Dataset

Regression with Two Datasets

Future Work

Regression Analysis of One Dataset



$$y = \begin{bmatrix} x \\ y \end{bmatrix} + \epsilon$$

Problem of Interest: which variables are associated with y?

Regression Analysis of One Dataset



$$y = \begin{bmatrix} x \\ y \end{bmatrix} + \epsilon$$

Problem of Interest: which variables are associated with y?

Common challenges:

- ightharpoonup number of variables p > n
- correlations among variables



$$\begin{split} \hat{\beta}^{\text{ridge}} &= \underset{\beta}{\text{arg min}} \left\{ \|y - X\beta\|^2 + \lambda \|\beta\|^2 \right\} \\ &= (X^{\mathsf{T}}X + \lambda I_p)^{-1}X^{\mathsf{T}}y \end{split}$$

The ridge estimator $\hat{eta}^{\mathrm{ridge}}$ has less variance, though it introduces bias.



$$\begin{split} \hat{\beta}^{\text{ridge}} &= \underset{\beta}{\text{arg min}} \left\{ \|y - X\beta\|^2 + \lambda \|\beta\|^2 \right\} \\ &= (X^{\mathsf{T}}X + \lambda I_p)^{-1}X^{\mathsf{T}}y \end{split}$$

The ridge estimator $\hat{eta}^{\mathrm{ridge}}$ has less variance, though it introduces bias.

Strongly positively correlated predictors have similar effects on y.



$$X = USV^{\mathsf{T}}$$

 $\hat{\beta}^{\mathrm{ridge}} = V(S^2 + \lambda I_n)^{-1}SU^{\mathsf{T}}y$



$$X = USV^{\mathsf{T}}$$

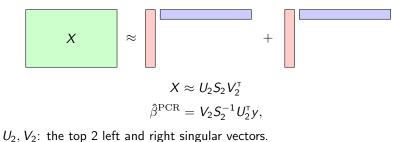
$$\hat{\beta}^{\mathrm{ridge}} = V(S^2 + \lambda I_n)^{-1}SU^{\mathsf{T}}y$$

Statistical inference is possible under certain conditions¹.

¹Bühlmann. Bernoulli, 13'

Principal Component Regression





Principal Component Regression



$$X pprox \mathcal{X} pprox \mathcal{X} pprox \mathcal{X} pprox \mathcal{U}_2 S_2 V_2^{ au} \ \hat{eta}^{ ext{PCR}} = V_2 S_2^{-1} U_2^{ au} y,$$

 U_2 , V_2 : the top 2 left and right singular vectors.

PCR projects high-dimensional data into 2D, and estimates association between the 2D summary and y. The PCR estimator $\hat{\beta}^{PCR}$ is effectively in 2D.

Latent Variable Model



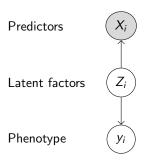


Figure: Illustration of the data generating mechanism underlying PCR. Both predictors and response are driven by the latent factors. PCR extracts latent factors via PCA of X and uses the resulting estimates \widehat{Z} as the new predictor. Regressing y against \widehat{Z} yields the PCR estimate, which typically enjoys smaller variance though it introduces bias.

²Wang et al. arXiv, 21'

Latent Variable Model



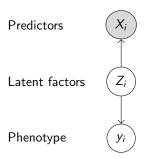


Figure: Illustration of the data generating mechanism underlying PCR. Both predictors and response are driven by the latent factors. PCR extracts latent factors via PCA of X and uses the resulting estimates \widehat{Z} as the new predictor. Regressing y against \widehat{Z} yields the PCR estimate, which typically enjoys smaller variance though it introduces bias.

Statistical inference is also possible under certain conditions².

²Wang et al. arXiv, 21'

L1 Regularization



L1 penalized regression projects HD data into a LD space defined by a few active variables, and estimates association between the LD summary and y.

The Lasso estimator $\hat{\beta}^{lasso}$ is effectively in LD.

Statistical inference is also possible under <u>compatibility</u> and <u>sparsity</u> conditions³.

³Zhang and Zhang, JRSSB, 14'; van de Geer et al., AoS, 14'; Cai and Guo, AoS, 17'

Naive Regression Analysis of Two Datasets

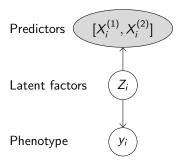


Could apply ridge regression or PCR by regressing y on $[X^{(1)}, X^{(2)}]$, but ...

Naive Regression Analysis of Two Datasets



Could apply ridge regression or PCR by regressing y on $[X^{(1)}, X^{(2)}]$, but ...



unable to tell what drives the latent factors: $Var(X^{(1)})$? $Var(X^{(2)})$? or $Cov(X^{(1)}, X^{(2)})$?

Outline



Motivation

Regression with One Datase

Regression with Two Datasets

Future Work

Regression Analysis of Two Datasets



Problem of Interest: which metabolites and bacteria are $\underline{jointly}$ associated with y?



CoRe⁴ performs supervised canonical correlation analysis by optimizing

$$b_1\|y - X^{(1)}\theta_1\|^2 + b_2\|y - X^{(2)}\theta_2\|^2 + b_3\|X^{(1)}\theta_1 - X^{(2)}\theta_2\|^2$$

⁴Gross and Tibshirani, Biostatistics, 15'



CoRe⁴ performs supervised canonical correlation analysis by optimizing

$$b_1\|y-X^{(1)}\theta_1\|^2+b_2\|y-X^{(2)}\theta_2\|^2+b_3\|X^{(1)}\theta_1-X^{(2)}\theta_2\|^2$$

CoRe uncovers signal that is common to $y, X^{(1)}$ and $X^{(2)}$.

⁴Gross and Tibshirani, Biostatistics, 15'



CoRe⁴ performs supervised canonical correlation analysis by optimizing

$$b_1\|y-X^{(1)}\theta_1\|^2+b_2\|y-X^{(2)}\theta_2\|^2+b_3\|X^{(1)}\theta_1-X^{(2)}\theta_2\|^2$$

CoRe uncovers signal that is common to $y, X^{(1)}$ and $X^{(2)}$.

► CoRe only uses the first canonical direction.

⁴Gross and Tibshirani, Biostatistics, 15'



 CoRe^4 performs supervised canonical correlation analysis by optimizing

$$b_1\|y-X^{(1)}\theta_1\|^2+b_2\|y-X^{(2)}\theta_2\|^2+b_3\|X^{(1)}\theta_1-X^{(2)}\theta_2\|^2$$

CoRe uncovers signal that is common to $y, X^{(1)}$ and $X^{(2)}$.

- ► CoRe only uses the first canonical direction.
- Statistical inference is not available for high-dimensional datasets.

⁴Gross and Tibshirani, Biostatistics, 15'

Canonical Variate Regression



CVR⁵ considers multiple phenotypes by optimizing

$$b_1 \| Y - X^{(1)} W_1 \theta_1 \|_F^2 + b_2 \| Y - X^{(2)} W_2 \theta_2 \|_F^2 + b_3 \| X^{(1)} W_1 - X^{(2)} W_2 \|_F^2$$

CVR uncovers signal that is common to $Y, X^{(1)}$ and $X^{(2)}$.

Canonical Variate Regression



CVR⁵ considers multiple phenotypes by optimizing

$$b_1 \| Y - X^{(1)} W_1 \theta_1 \|_F^2 + b_2 \| Y - X^{(2)} W_2 \theta_2 \|_F^2 + b_3 \| X^{(1)} W_1 - X^{(2)} W_2 \|_F^2$$

CVR uncovers signal that is common to $Y, X^{(1)}$ and $X^{(2)}$.

► Statistical inference is not available for high-dimensional datasets.

Latent Variable Model



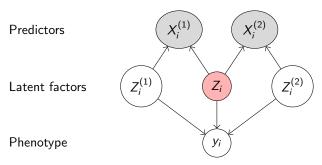


Figure: Illustration of the data generating mechanism underlying JMDR. Two sets of predictors are driven by latent factors that are shared (Z) and dataset-specific $(Z^{(d)})$. All latent factors can contribute to y, though we are mainly interested in the shared information.



<u>Our Hypothesis:</u> Information shared between datasets may be more causal.



Our Hypothesis: Information shared between datasets may be more causal.

Shared factors Z_i capture joint variation that is more likely due to the underlying biological mechanism.



<u>Our Hypothesis:</u> Information shared between datasets may be more causal.

Shared factors Z_i capture joint variation that is more likely due to the underlying biological mechanism.

Individual factors $Z_i^{(1)}$ and $Z_i^{(2)}$ capture dataset-specific variation due to, e.g. batch effects, technical artifacts, host factors, etc.



<u>Our Hypothesis:</u> Information shared between datasets may be more causal.

Shared factors Z_i capture joint variation that is more likely due to the underlying biological mechanism.

Individual factors $Z_i^{(1)}$ and $Z_i^{(2)}$ capture dataset-specific variation due to, e.g. batch effects, technical artifacts, host factors, etc.

Separating individual from joint factors may improve estimation.



Add a slide on biological interpretations of the latent factors





- \blacktriangleright W_1 and W_2 are low-rank loading matrices for *joint variation*.
- V₁ and V₂ are low-rank loading matrices for dataset-specific variation.



$$X^{(1)} = \mathbf{Z} W_1^{\mathsf{T}} + Z^{(1)} V_1^{\mathsf{T}} + E^{(1)}$$

$$X^{(2)} = \mathbf{Z} W_2^{\mathsf{T}} + Z^{(2)} V_2^{\mathsf{T}} + E^{(2)}$$

$$y = X^{(1)} \beta_1^* + X^{(2)} \beta_2^* + \epsilon$$

- \blacktriangleright W_1 and W_2 are low-rank loading matrices for *joint variation*.
- V₁ and V₂ are low-rank loading matrices for dataset-specific variation.
- Latent factors are independent.
- ▶ The errors ϵ and $E^{(d)}$ have mean zero, and are independent of each other.



$$X^{(1)} = \mathbf{Z} W_1^{\mathsf{T}} + Z^{(1)} V_1^{\mathsf{T}} + E^{(1)}$$

$$X^{(2)} = \mathbf{Z} W_2^{\mathsf{T}} + Z^{(2)} V_2^{\mathsf{T}} + E^{(2)}$$

$$y = X^{(1)} \beta_1^* + X^{(2)} \beta_2^* + \epsilon$$

- ▶ W_1 and W_2 are low-rank loading matrices for *joint variation*.
- V₁ and V₂ are low-rank loading matrices for dataset-specific variation.
- Latent factors are independent.
- ▶ The errors ϵ and $E^{(d)}$ have mean zero, and are independent of each other.
- ▶ Not interested in β_1^* and β_2^* → Need to define target parameter



Letting
$$W^{\mathsf{T}} = [W_1^{\mathsf{T}}, W_2^{\mathsf{T}}],$$

$$y pprox \left(ZW^{\mathsf{T}} + Z^{(1)}V_1^{\mathsf{T}} + Z^{(2)}V_2^{\mathsf{T}} \right) egin{pmatrix} eta_1^* \ eta_2^* \end{pmatrix}$$

► The quantity

$$\alpha = W^{\mathsf{T}} \binom{\beta_1^*}{\beta_2^*} \in \mathbb{R}^k$$

reflects the association between the joint variation and y.



Letting
$$W^{\scriptscriptstyle\mathsf{T}} = [W_1^{\scriptscriptstyle\mathsf{T}}, W_2^{\scriptscriptstyle\mathsf{T}}],$$

$$y pprox \left(\mathbf{Z} W^{\mathsf{T}} + Z^{(1)} V_1^{\mathsf{T}} + Z^{(2)} V_2^{\mathsf{T}} \right) \begin{pmatrix} \beta_1^* \\ \beta_2^* \end{pmatrix}$$

Each coordinate in

$$heta = WW^{\mathsf{T}}inom{eta_1^*}{eta_2^*} \in \mathbb{R}^{
ho_1+
ho_2}$$

reflects the association between each variable and y through joint variation.

We use θ to determine which metabolites and bacteria are <u>jointly</u> associated with y.

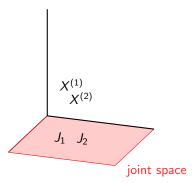
Model Identifiability



- $ightharpoonup \alpha$ is identifiable up to rotations
- \triangleright θ is invariant to rotations

Understanding Joint Association





 $X^{(1)} pprox J_1, X^{(2)} pprox J_2$ in the joint space, but J_1 and J_2 are special low-rank matrices

 θ is effectively low-dimensional

Comparison to Linear Regulatory Module



LRM⁶ considers the relationship $X_i^{(1)} o X_i^{(2)} o y_i$

$$X^{(2)}W_2 = X^{(1)}W_1 + E$$

 $y = X^{(2)}W_2\beta_0 + res_1 \cdot \beta_1 + res_2 \cdot \beta_2 + \epsilon$

Columns in W_1 and W_2 are assumed to be sparse, reflecting a small number of active regulation.

Comparison to Linear Regulatory Module



LRM⁶ considers the relationship $X_i^{(1)} o X_i^{(2)} o y_i$

$$X^{(2)}W_2 = X^{(1)}W_1 + E$$

$$y = X^{(2)}W_2\beta_0 + res_1 \cdot \beta_1 + res_2 \cdot \beta_2 + \epsilon$$

Columns in W_1 and W_2 are assumed to be sparse, reflecting a small number of active regulation.

LRM does not provide statistical inference for selected variables.

⁶Zhu et al. Biostatistics, 16'



How JMDR works (briefly)



$$\chi^{(1)}$$
 = $+\cdots+$

$$U_1 = U_2 = U_2 = U_3$$
 $n \text{ sample scores from } X^{(1)}$
 $n \text{ sample scores from } X^{(2)}$



We use $U_1^T U_2$ to project $X^{(1)}$ and $X^{(2)}$ onto a *joint subspace*.



We use $U_1^T U_2$ to project $X^{(1)}$ and $X^{(2)}$ onto a joint subspace.

▶ Decompose $U_1^T U_2 = RDQ^T$ and define a projection

$$J_1 = (U_1 R R^{\mathsf{T}} U_1^{\mathsf{T}}) X^{(1)}$$

$$J_2 = (U_2 Q Q^{\mathsf{T}} U_2^{\mathsf{T}}) X^{(2)}$$

$$J_2 = (U_2 Q Q^{\mathsf{T}} U_2^{\mathsf{T}}) X^{(2)}$$



We use $U_1^T U_2$ to project $X^{(1)}$ and $X^{(2)}$ onto a *joint subspace*.

▶ Decompose $U_1^{\mathsf{T}}U_2 = RDQ^{\mathsf{T}}$ and define a projection

$$J_1 = (U_1 R R^{\mathsf{T}} U_1^{\mathsf{T}}) X^{(1)}$$

$$J_2 = (U_2 Q Q^{\mathsf{T}} U_2^{\mathsf{T}}) X^{(2)}$$

Remark: Use principal angles to determine the dimensions of $U_1^T U_2$ that can be considered *joint*.



We use $U_1^T U_2$ to project $X^{(1)}$ and $X^{(2)}$ onto a *joint subspace*.

▶ Decompose $U_1^{\mathsf{T}}U_2 = RDQ^{\mathsf{T}}$ and define a projection

$$J_1 = (U_1 R R^{\mathsf{T}} U_1^{\mathsf{T}}) X^{(1)}$$

$$J_2 = (U_2 Q Q^{\mathsf{T}} U_2^{\mathsf{T}}) X^{(2)}$$

<u>Remark</u>: Use principal angles to determine the dimensions of $U_1^T U_2$ that can be considered *joint*.

ightharpoonup Estimate the latent factors by $[J_1,J_2]=\widehat{Z}\widehat{\Lambda}\widehat{W}^{\scriptscriptstyle\mathsf{T}}.$



We use $U_1^T U_2$ to project $X^{(1)}$ and $X^{(2)}$ onto a *joint subspace*.

▶ Decompose $U_1^{\mathsf{T}}U_2 = RDQ^{\mathsf{T}}$ and define a projection

$$J_1 = (U_1 R R^{\mathsf{T}} U_1^{\mathsf{T}}) X^{(1)}$$

$$J_2 = (U_2 Q Q^{\mathsf{T}} U_2^{\mathsf{T}}) X^{(2)}$$

<u>Remark</u>: Use principal angles to determine the dimensions of $U_1^T U_2$ that can be considered *joint*.

- ightharpoonup Estimate the latent factors by $[J_1,J_2]=\widehat{Z}\widehat{\Lambda}\widehat{W}^{\scriptscriptstyle\mathsf{T}}.$
- ▶ Regress y on \widehat{Z} .

Statistical Inference



Global test:

$$H_0: \alpha = 0$$

Element-wise test:

$$H_{0,j}:\theta_j=0$$

JMDR inference accounts for subspace recovery error.



Results

BactoCARB Study



```
X^{(1)} = 144 metabolites & bile acids (blood)
```

 $X^{(2)} = 134$ microbial genus abundances (stool)

 $y = {\sf enterolactone}$ 'ENL' (urine)

BactoCARB Study



 $X^{(1)} = 144$ metabolites & bile acids (blood)

 $X^{(2)} = 134$ microbial genus abundances (stool)

y =enterolactone 'ENL' (urine)

Rank of joint space: 4

Rank of individual space: 11 for metabolites and 4 for microbiome

BactoCARB Study



 $X^{(1)} = 144$ metabolites & bile acids (blood)

 $X^{(2)} = 134$ microbial genus abundances (stool)

y =enterolactone 'ENL' (urine)

Rank of joint space: 4

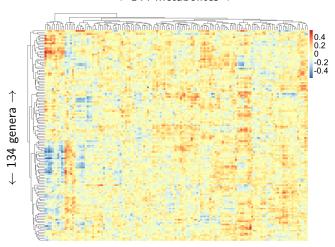
Rank of individual space: 11 for metabolites and 4 for microbiome

A global test for association between ENL and \underline{joint} signal: $p\text{-value} < 10^{-5}$

Correlations btw Metabolites and Microbes



\leftarrow 144 metabolites \rightarrow



Correlations in the Joint Space



\leftarrow 144 metabolites \rightarrow

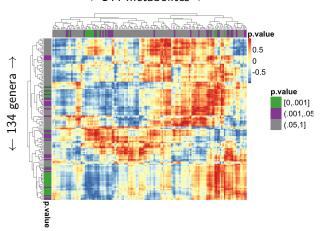
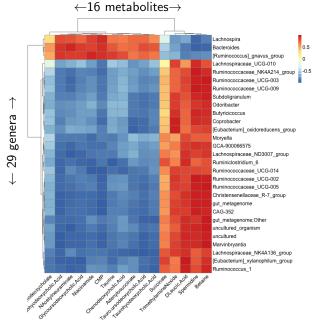


Figure: Pairwise correlations in the 4-dim joint space

p-values are from element-wise test in JMDR model ENL \sim metabolites + genera

Significant Variables





Summary



► Given 2 (or more) datasets, strip off individual/separate variation and represent data using only shared/joint variation.

Summary



- ► Given 2 (or more) datasets, strip off individual/separate variation and represent data using only shared/joint variation.
- ► Joint Matrix Decomposition Regression identifies variables in each dataset that are jointly associated with a phenotype.

Summary



- ► Given 2 (or more) datasets, strip off individual/separate variation and represent data using only shared/joint variation.
- ► Joint Matrix Decomposition Regression identifies variables in each dataset that are jointly associated with a phenotype.
- ▶ JMDR also works for dichotomous outcomes.

Dog Aging Project



An epidemiological cohort $n \ge 1000$ of companion dogs aimed at identifying the genetic and environmental factors of aging

- $ightharpoonup X^{(1)} = \text{metabolites (blood)}$
- $ightharpoonup X^{(2)} = microbial species/gene abundances (stool)$
- $ightharpoonup X^{(3)} = \text{epigenetic markers (PBMCs)}$
- ightharpoonup y = cognitive score, aging score

Problem of Interest: which metabolites, bacteria, and epigenetic markers are *jointly* associated with *y*?

Future Work



- Extend to more than 2 datasets and partially shared joint signals
- Extend from linear to nonlinear models
- Bayesian modeling
- Applications to uncovering biomarkers related to cognition and aging (Dog Aging Project, HCHS/SOL)

Aknowledgement









Ali Shojaie @UW



Tim Randolph @Fred Hutch

Thank You!

https://drjingma.com