

Morning session: overview

Unsupervised learning: dimension reduction and data integration

- Principal Component Analysis (PCA)
 - Maximal variance principle
- Partial Least Squares (PLS)
 - Maximal covariance principle
- Two-way Orthogonal PLS (O2PLS)
 - Population model
 - Multi-omics data integration



Introduction

- We have seen PCA and PLS
- Looked at them from an algorithmic point of view
- What is the underlying population model?
 - What are the random variables, and what are the parameters?



Population model for PCA

$$x = tW^{\top} + e$$

- t are scores (the PCs)
- W are weights (or loadings), parameters
- e are residuals
- Parameters are estimated such that the variance of Xw is maximized



Population model for PLS

$$x = tW^{\top} + e$$

$$y = uC^{\top} + f$$

$$u = tB + h$$

- t and u are latent (hidden) scores (the joint PCs)
- w and c are weights (or loadings), parameters
- The relation between x and y is fully captured by t and u , via the third relation
- Parameters are estimated such that the covariance between Xw and Yc is maximized



Interpretation of the model

$$\begin{aligned}x &= tW^{\top} + e \\y &= uC^{\top} + f \\u &= tB + h\end{aligned}$$

- x and y are two datasets, say genetic and metabolomic data
- t and u are latent variables underlying these data
- These latent variables vary, and through W and C cause variation in the two datasets
- In the above context, t and u could be genetic/metabolic pathways of several connected genes and metabolites
- Then, W and C would tell us which genes and metabolites are involved



Omics-specific variation

- PCA models the variance of X (or Y)
- PLS models the covariance
- Suppose both sources are present, independently
 - E.g. some pathways connect genes and metabolites
 - Other pathways are there for self-maintenance
- Capture both parts at the same time
- Need to extend the PLS model



O2PLS model and data-specific parts

$$\begin{aligned}x &= tW^{\top} + t_s W_s^{\top} + e \\y &= uC^{\top} + u_s C_s^{\top} + f \\u &= tB + h\end{aligned}$$

- x and y are two datasets, say genetic and metabolomic data
- In addition to the PLS joint weights and scores, we have specific weights and scores
- The relation between x and y is still fully captured by t and u



Estimating the O2PLS components

- Three step estimation
 - First estimate several PLS components, they will contain both joint and specific parts
 - From that, estimate specific parts only and subtract
 - Finally, estimate again PLS on the “corrected” data
- Software available from CRAN: *OmicsPLS*



Number of components

- Until now, we did not mention how to find out the number of components needed
- For PCA, this number is the number of PCs
- For PLS, this is the number of joint PCs
- For O2PLS, this is
 - The number of Joint PCs
 - And the number of X-specific components
 - And Y-specific components
- Standard way is to do cross-validation (next session)
 - Pick a number, then evaluate the performance on an independent test set
 - Repeat, and select the one that gives the best performance



O2PLS: summary

- For given datasets X and Y , we want to inspect their relation
- Want to capture data-specific parts as well
- O2PLS estimates joint and specific parts, consisting of weights and scores
- One can interpret or plot the weights and scores to understand which features/samples are most important



Morning session: overview

Unsupervised learning: dimension reduction and data integration

- Principal Component Analysis (PCA)
 - Maximal variance principle
- Partial Least Squares (PLS)
 - Maximal covariance principle
- Two-way Orthogonal PLS (O2PLS)
 - Population model
 - Multi-omics data integration

