

ETC3250/5250: Regularization

Semester 1, 2020

Professor Di Cook

Econometrics and Business Statistics

Monash University

Week 9 (a)

Too many variables

Fitting a linear regression model requires:

$$\begin{aligned} & \underset{\beta \in \mathbb{R}^p}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \\ & \equiv \underset{\beta \in \mathbb{R}^p}{\text{minimize}} (y - X\beta)'(y - X\beta) \end{aligned}$$

The least square solution for β is

$$\hat{\beta} = (X'X)^{-1}X'y$$

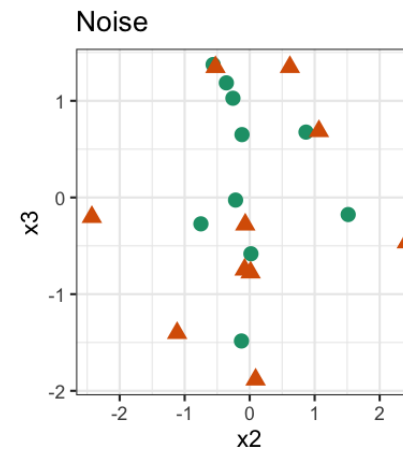
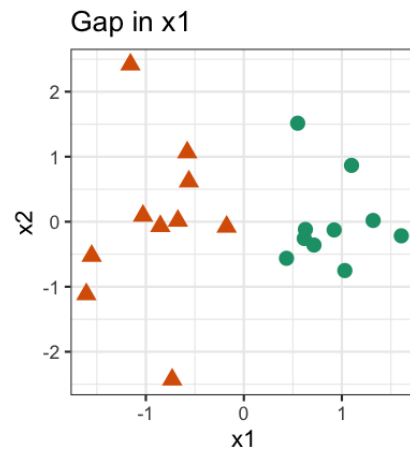
To **invert** a matrix, requires it to be **full rank**.

Example: Using simulation

20 observations

2 classes: A, B

One variable with separation, 99 noise variables



00:23

What will be the optimal LDA coefficients?

Fit linear discriminant analysis on first two variables.

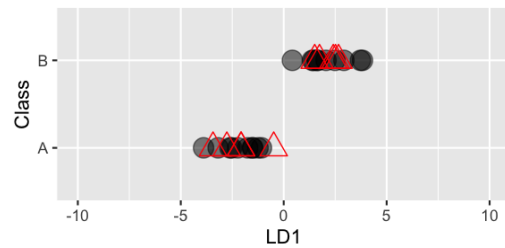
```
## Call:
## lda(cl ~ ., data = tr[, c(1:2, 101)], prior = c(0.5, 0.5))
##
## Prior probabilities of groups:
##   A   B
## 0.5 0.5
##
## Group means:
##           x1           x2
## A  0.8918346  0.0009586256
## B -0.8918346 -0.0009586256
##
## Coefficients of linear discriminants:
##           LD1
## x1 -2.41606038
## x2  0.05224863
```

Coefficient for **x1** MUCH higher than **x2**. As expected!

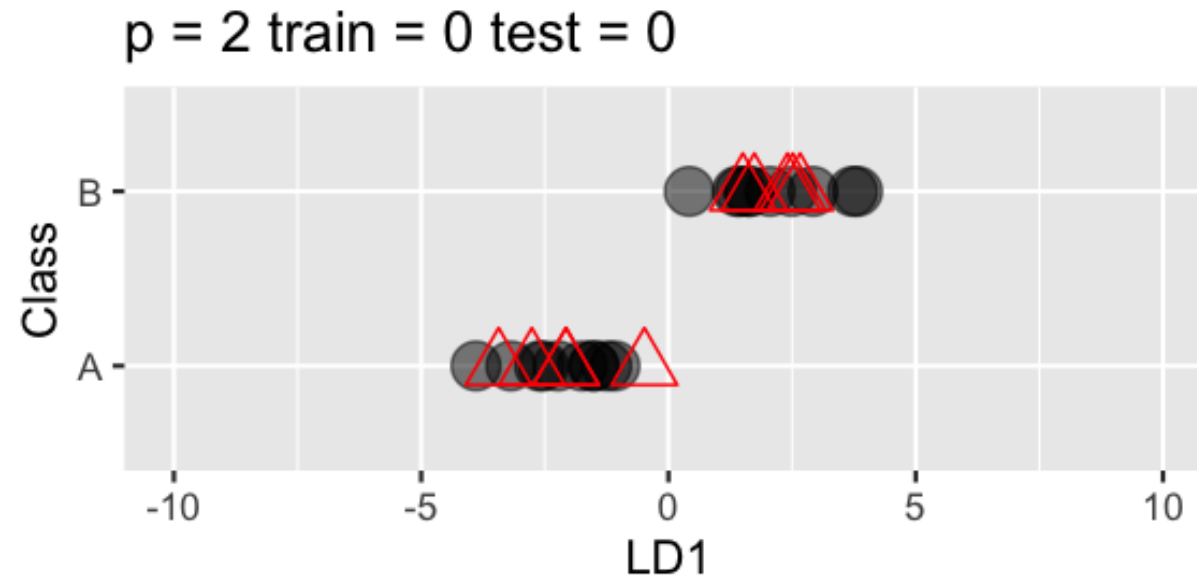
Predict the training and test sets

```
##  
##      A  B  
## A 10  0  
## B  0 10
```

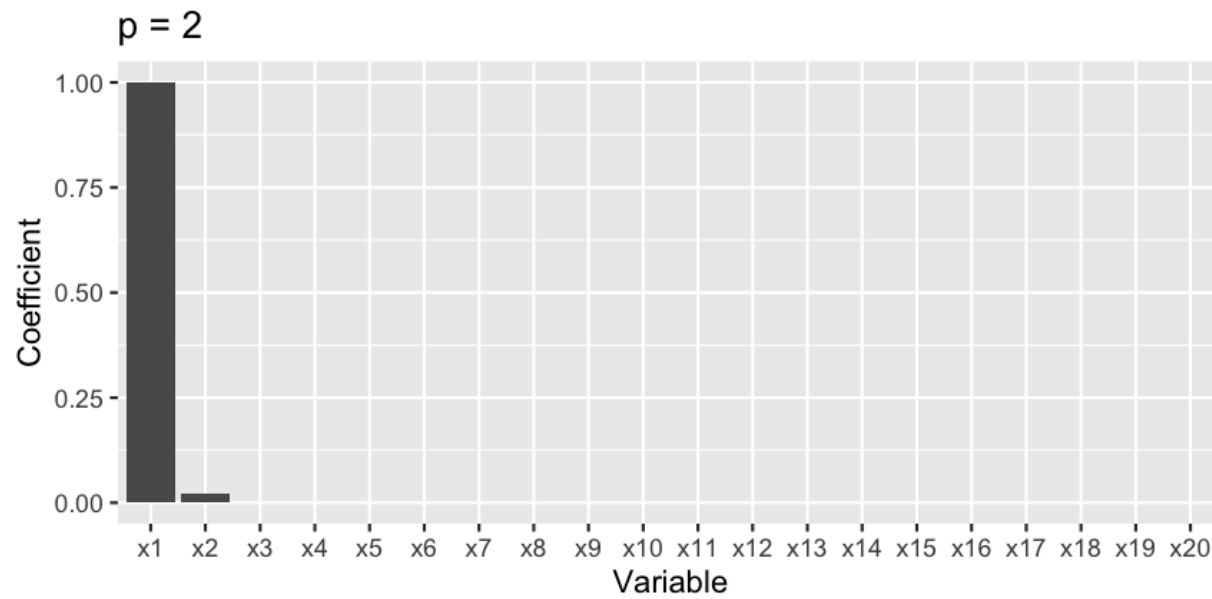
```
##  
##      A B  
## A 5 0  
## B 0 5
```



What happens to test set (and predicted training values) as number of noise variables increases:



Estimated coefficients as dimensions of noise increase:



How do we tackle high-dimension, low sample size problems?

Subset selection

Identify a subset s of the p predictors, most related to response.

$$\begin{aligned} & \underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \\ & \text{subject to } \sum_{j=1}^p I(\beta_j \neq 0) \leq k, \quad k \geq 0. \end{aligned}$$

where $k \geq 0$ is a tuning parameter.

- Need to consider $\binom{p}{k}$ models containing s predictors
computationally infeasible when p and s are large
- Stepwise procedures: forward, backward, etc.

Model fit statistics

These can be used to decide on choice of k .

▮ $MSE = RSS/n$, but the training MSE is an under-estimate of test MSE , and it will decrease with larger p .

▮ Methods for adjusting the training error for model size include Mallows C_p , Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC) and adjusted R^2 .

Mallows C_p

For a fitted least squares model containing d predictors, a reasonable estimate of the test MSE is:

$$C_p = \frac{1}{n}(RSS + 2d\hat{\sigma}^2)$$

where $\hat{\sigma}^2$ is an estimate of the variance of the error ε , computed from the full model containing all predictors.

The additional part penalises the training RSS to adjust for the under-estimation of test error.

AIC and BIC

$$AIC = \frac{1}{n\hat{\sigma}^2}(RSS + 2d\hat{\sigma}^2)$$

and hence is $\propto C_p$.

$$BIC = \frac{1}{n\hat{\sigma}^2}(RSS + \log(n)d\hat{\sigma}^2)$$

all tend to take on low values for models with small test error.

Adjusted R^2

$$\text{Adjusted } R^2 = 1 - \frac{RSS/(n - d - 1)}{TSS/(n - 1)}$$

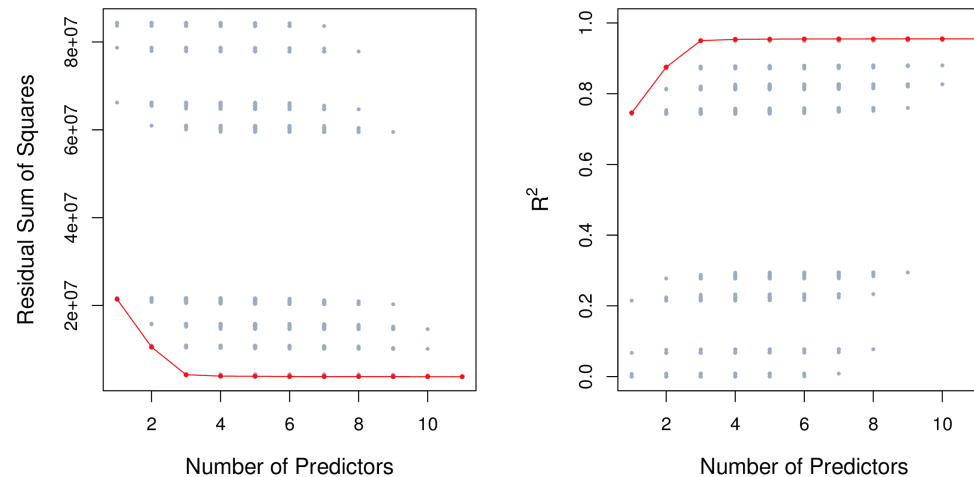
The intuition is that once all of the correct variables have been included in the model, adding additional *noise* variables will lead to only a very small decrease in RSS.

Best subset selection algorithm

1. Let \mathcal{M}_o denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 1, 2, \dots, p$:
 - a. Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - b. Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Best means smallest RSS (or largest R^2).
3. Select a single best model from among $\mathcal{M}_o, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .

Best subset selection algorithm

Best subset selection algorithm applied to the 11 predictors of the Credit data.



Forward stepwise selection

Forward stepwise selection is a computationally efficient alternative to best subset selection. It considers a much smaller set of models.

When $p = 20$, best subset selection requires fitting 1,048,576 models, whereas forward stepwise selection requires fitting only 211 models.

Forward stepwise selection - algorithm

1. Let \mathcal{M}_0 denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 0, 1, 2, \dots, p - 1$:
 - a. Consider all $p - k$ models that augment \mathcal{M}_k with *one additional predictor*.
 - b. Pick the best among these $p - k$ models, and call it \mathcal{M}_{k+1} . Best means smallest RSS (or largest R^2).
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .

Backwise stepwise selection

- Backward stepwise starts with all variables in the model, and removes the variable with smallest RSS.
- Forward and backwards stepwise procedures are not guaranteed to provide the best model.
- Backwards stepwise requires that $n > p$, but forward stepwise does not, and can stop adding variables once $n(< p)$ is reached.

Shrinkage methods

Shrinkage methods fit a model containing all p predictors using a technique that constrains or regularizes the coefficient estimates, or equivalently, that shrinks some of the coefficient estimates towards zero.

There are two main methods: Ridge regression and Lasso.

Ridge regression

$$\text{RSS} = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

Least squares:

$$\underset{\beta}{\text{minimize}} \text{RSS}$$

Ridge regression:

$$\underset{\beta}{\text{minimize}} \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$$

where $\lambda \geq 0$ is a tuning parameter.

Ridge regression

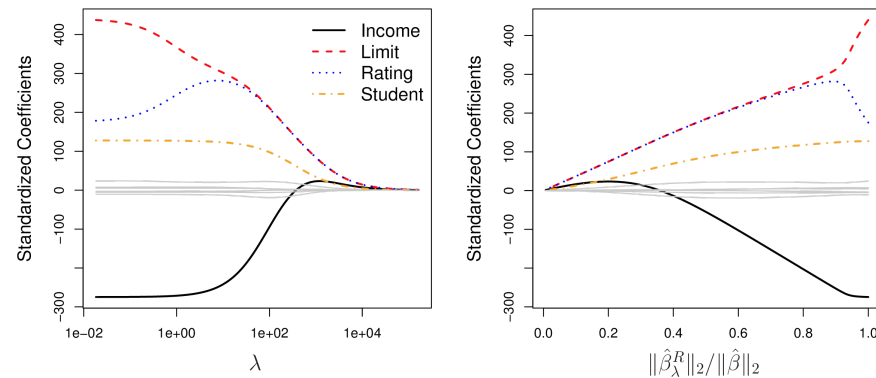
$$\lambda \sum_{j=1}^p \beta_j^2$$

is called a **shrinkage penalty**. It is small when β_1, \dots, β_p are close to 0.


λ serves as a **tuning parameter**, controlling the relative impact of these two terms on the regression coefficient estimates. When it is 0, the penalty term has no effect on the fit.

Ridge regression will produce a **different set of coefficients** for each λ , call them $\hat{\beta}_{\lambda}^R$. Tuning λ , typically by cross-validation, is critical component of fitting the model.


Standardized ridge regression coefficients for the Credit data set.





(Chapter6/6.4.pdf)

 $p = 10$

 Left side of plot corresponds to least squares.

 When λ is extremely large, then all of the ridge coefficient estimates are basically zero, which is the null model.

 4 of 10 variables have larger coefficients, and one, Rating, initially increases with λ .

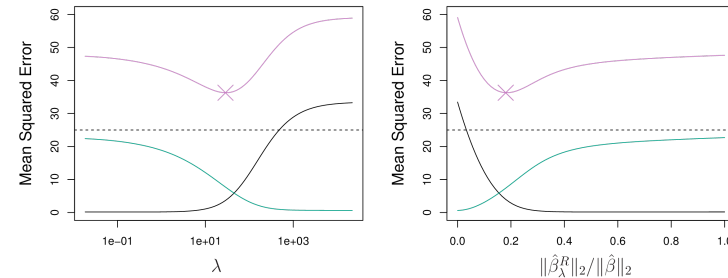
 Right-side plot, x -axis indicates amount the coefficients shrink to 0, value of 1 indicates LS.

The scale of variables can affect ridge regression performance.

It is important to standardise the scale of predictors prior to ridge regression.

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sigma_{x_j}}$$

Simulation scenario! Ridge regression improves on least squares, for large number of variables, in the bias-variance tradeoff. It sacrifices some bias for the benefit of decreased variance.



bias variance test error

The Lasso

Ridge regression:

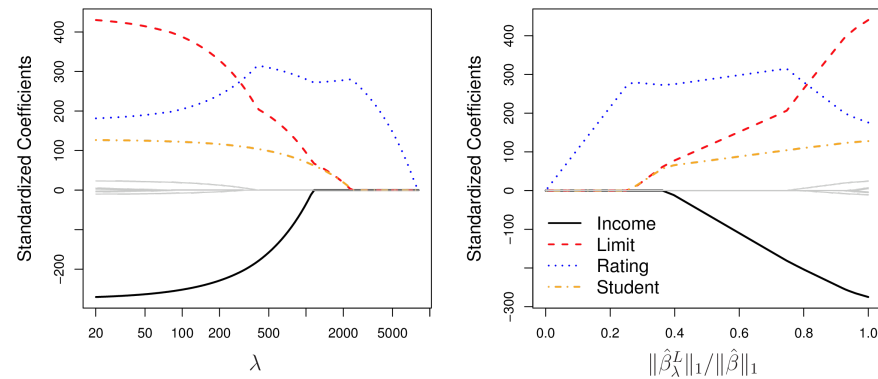
$$\underset{\beta}{\text{minimize}} \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$$

Lasso:

$$\underset{\beta}{\text{minimize}} \text{RSS} + \lambda \sum_{j=1}^p |\beta_j|$$

and same $\lambda \geq 0$ is a tuning parameter.

Standardized lasso coefficients for the Credit data set.



(Chapter6/6.6.pdf)

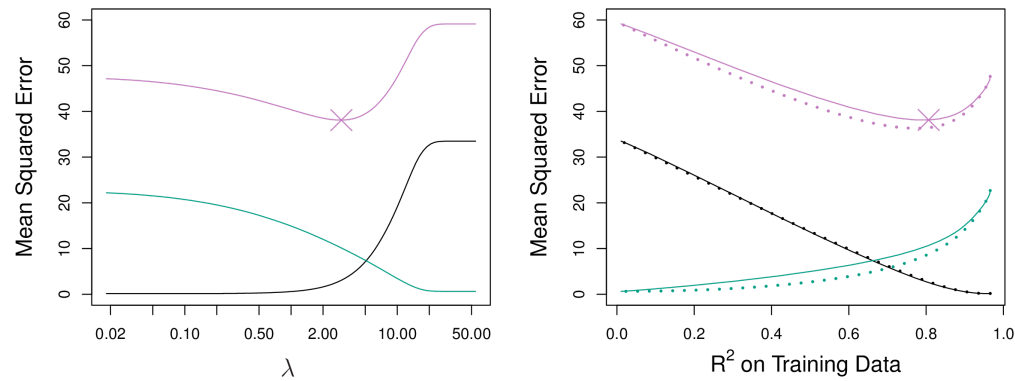
|| $p = 10$

|| Has the effect of forcing some variables exactly to 0.

|| Cleaner solution than ridge regression.

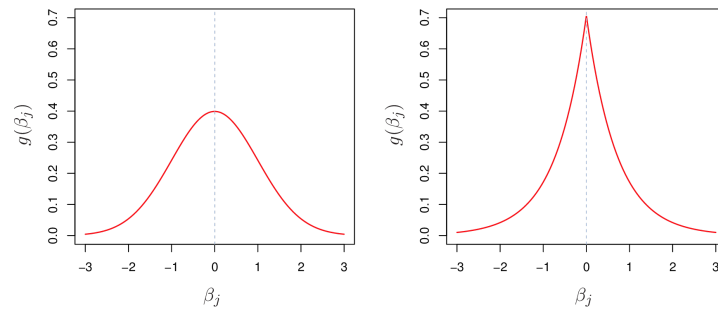
Simulation scenario!

Bias-variance tradeoff with lasso, and comparison against ridge regression.



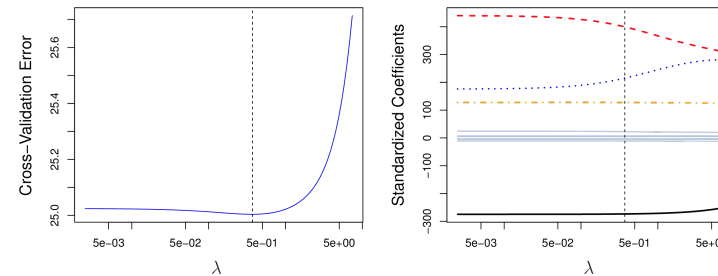
Bias Variance Test error

Bayesian interpretation: Ridge regression is the posterior mode for β under a Gaussian prior (left); The lasso is the posterior mode for β under a double-exponential prior (right).



(Chapter6/6.11.pdf)

Cross-validation on the Credit example, yields a suggestion to use $\lambda = 0.5$ for ridge regression model.



(Chapter6/6.12.pdf)

Penalised LDA

Recall: LDA involves the eigen decomposition of $\Sigma^{-1}\Sigma_B$, where

$$\Sigma_B = \frac{1}{K} \sum_{i=1}^K (\mu_i - \mu)(\mu_i - \mu)'$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_i)(x_i - \mu_i)'$$

The eigen-decomposition is an analytical solution to a sequential optimisation problem:

$$\underset{\beta_k}{\text{maximize}} \beta_k^T \hat{\Sigma}_B \beta_k$$

$$\text{subject to } \beta_k^T \hat{\Sigma}_B \beta_k \leq 1, \beta_k^T \hat{\Sigma}_B \beta_j = 0 \quad \forall i < k$$

Penalised LDA

The problem is inverting Σ^{-1} , fix it by

$$\begin{aligned} & \underset{\beta_k}{\text{maximize}} \left(\beta_k^T \hat{\Sigma}_B \beta_k + \lambda_k \sum_{j=1}^p |\hat{\sigma}_j \beta_{kj}| \right) \\ & \text{subject to } \beta_k^T \tilde{\Sigma} \beta_k \leq 1 \end{aligned}$$

where $\hat{\sigma}_j$ is the within-class standard deviation for variable j . This is **penalised LDA**, and see [reference](#), and the [R package](#).

PDA Index

Penalised LDA projection pursuit index. Available in the `tourr` package.

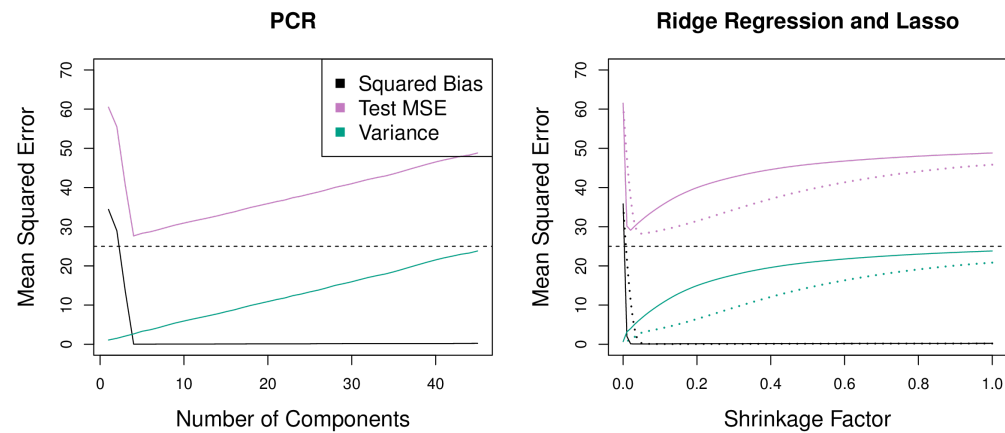
$$I_{PDA}(A, \lambda) = 1 - \frac{\left| A' \{ (1 - \lambda) \hat{\Sigma} + n \lambda I_p \} A \right|}{\left| A' \{ (1 - \lambda) (\hat{\Sigma}_B + \hat{\Sigma}) + n \lambda I_p \} A \right|}$$

Optimising this function over $p \times d$ projection matrix A .

Principal component regression

The principal components regression (PCR) approach involves constructing the first M principal components, Z_1, \dots, Z_M , and then using these components as the predictors in a linear regression model, that is fit using least squares.

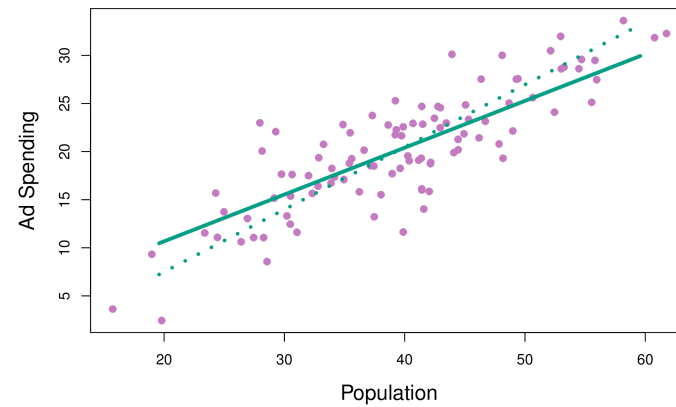
PCR, ridge regression, and the lasso compared on simulated data. PCR does well when the response is related to few PCs.



Bias Variance Test error

Partial least squares

Partial least squares (PLS), a supervised alternative to PCR.



Two predictors are shown: Solid line is PLS, dashed line is PCR.

Partial least squares

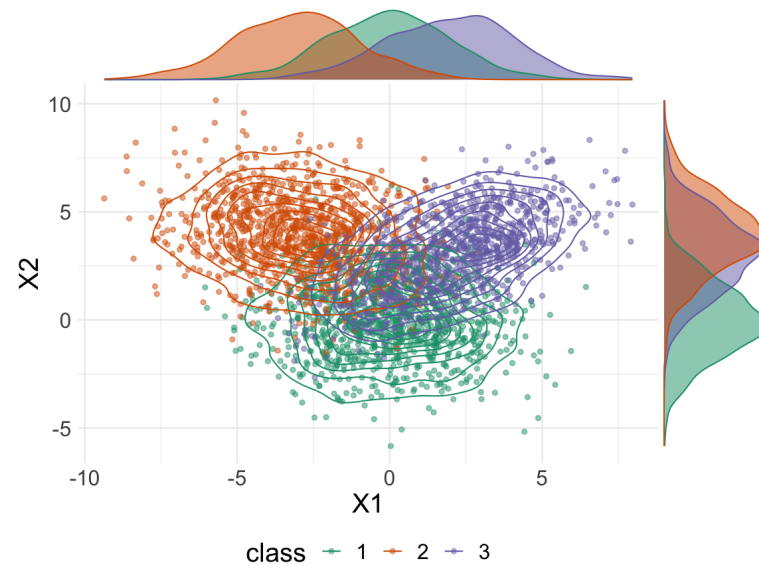
1. Standardise all variables
2. Find $Z_1 = \phi_{1j}X_j$ by setting ϕ_{1j} to be the coefficient from a simple linear regression model $Y \sim X_j$.
3. To find Z_2 , first regress each variable on Z_1 and use the residuals, call these X_j^r . Then find $Z_2 = \phi_{2j}X_j^r$ by setting ϕ_{2j} to be the coefficient from a simple linear regression model $Y \sim X_j^r$.
4. Repeat steps 2-3 until we have Z_1, \dots, Z_M .

Final model fitted for Y using Z_1, \dots, Z_M .

Performance is no better than ridge regression or PCR. Can reduce bias, has potential to increase variance. PLS is similar to partial regression, where new variables are first regressed on predictors that are already in the model, and it is the residuals that are used.

Diagonal Discriminant Analysis

- The simplest form of regularisation assumes that the features are independent within each class.
- Consider a *diagonal-covariance* LDA rule for classifying classes
- A special case of the naive-Bayes classifier



Discriminant Function

It can be shown that the discriminant score for a new observation \mathbf{x}^* when the features are considered independent reduces to the following:

$$\delta_k(\mathbf{x}^*) = - \sum_{j=1}^p \frac{(x_j^* - \bar{x}_{kj})^2}{s_j^2} + 2 \log \pi_k.$$

The classification rule is then

$$C(\mathbf{x}^*) = \ell \quad \text{if} \quad \delta_\ell(\mathbf{x}^*) = \max_k \delta_k(\mathbf{x}^*).$$

Filter features for prediction

To motivate the upcoming method, consider a binary classification DLDA problem.

One way we could establish which of the features are driving prediction would be to perform a two-sample t -test

$$t_j = \frac{\bar{x}_{1j} - \bar{x}_{0j}}{s_j}$$

with the t statistic providing a measure of how significant the difference in class means for predictor j .

Filter features for prediction

Think about it: Using the t statistic - $t_j = \frac{\bar{x}_{1j} - \bar{x}_{0j}}{s_j}$ for all features, what is one way we can determine important features for prediction?

Filter features for prediction

Answer: Can consider filtering for features with $|t_j| > 2$, as this is deemed significant at the 5% level.



Note - further consideration can be given to the issue of *Multiple Testing*

Nearest Shrunk Centroids (NSC)

Now consider the following statistic,

$$d_{kj} = \frac{\bar{x}_{kj} - \bar{x}_j}{m_k(s_j + s_0)} \quad \text{with} \quad m_k^2 = \frac{1}{N_k} - \frac{1}{N}$$

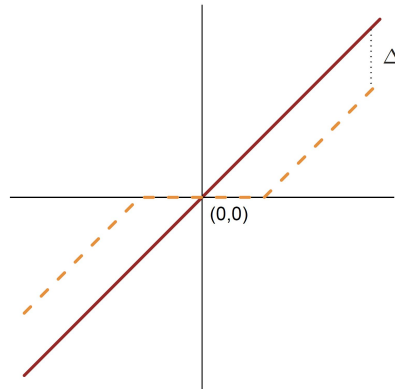
and s_0 a small value to protect d_{kj} from small expression values.

This statistic is a measure for how significant the difference between the class k mean for predictor j , and the overall mean for predictor j .

Soft Thresholding

Each d_{kj} is reduced by an amount Δ in absolute value, and is set to zero if its absolute value is less than zero.

$$d'_{kj} = \text{sign}(d_{kj})(|d_{kj}| - \Delta)_+,$$



Nearest Shrunk Centroids Classifier

The NSC uses either version of the statistic d'_{kj} to regularise by shrinking the class means towards the overall mean for each predictor separately as follows:

$$\bar{x}'_{kj} = \bar{x}_j + m_k(s_j + s_0)d'_{kj}$$

Unless a predictor has a significant difference to the overall mean for at least one class, it is useless for classification.

We then use the shrunk centroids \bar{x}'_{kj} in place of \bar{x}_{kj} in the DLDA discriminant function.

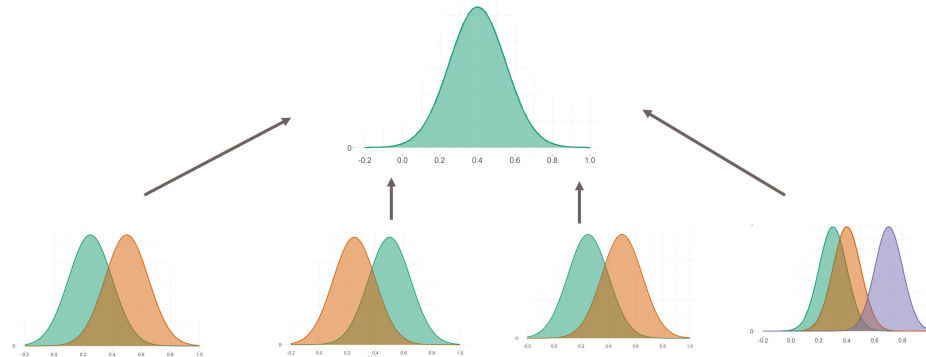
Alternative - penalised multiple hypothesis testing (multiDA)

Another approach to high dimensional DA involves formulating the problem as a multiple hypothesis testing problem, and asking the question - "What defines a discriminative feature?", and then choosing discriminative features through a penalised likelihood ratio test.

LRT – compare to the null

For $K = 3$ classes, there are $m = 5$ potential partitions of the data.

For all 5 hypotheses, compare the likelihood to the null. Pick the "partition" that is the most likely.



A penalised likelihood ratio test statistic

Two forms of penalisation can be considered:

 **The BIC** - useful when Positive Selection Rate is preferred to controlling False Discovery Rate (FDR).

$$\nu_m \log(n)$$

 **The Extended BIC** - useful for high dimensional data, penalising additionally on the number of features p .

$$\nu_m [\log(n) + 2 \log(p)]$$

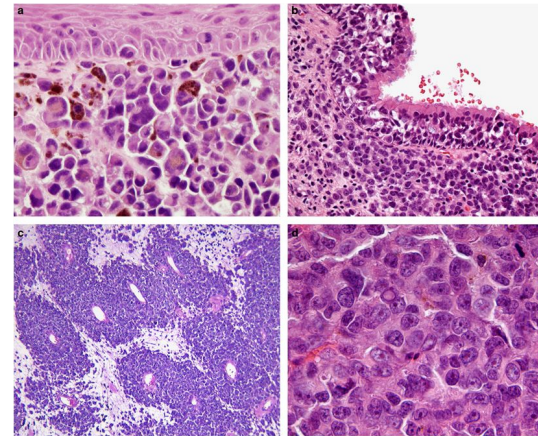
(Note - $\nu_m = g_m - 1$ where g_m is the number of groupings considered in model m).

Recall - SRBCT cancer prediction

▮ The SRBCT dataset (Khan et al., 2001) looks at classifying 4 classes of different childhood tumours sharing similar visual features during routine histology.

▮ Data contains 83 microarray samples with 1586 features.

▮ We will revisit this data later on in the course to explore high dimensional DA. Now is that time



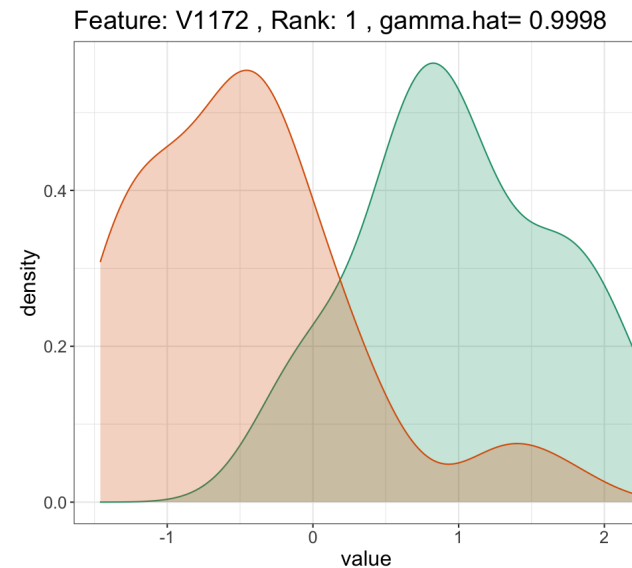
Source: Nature

multiDA in R

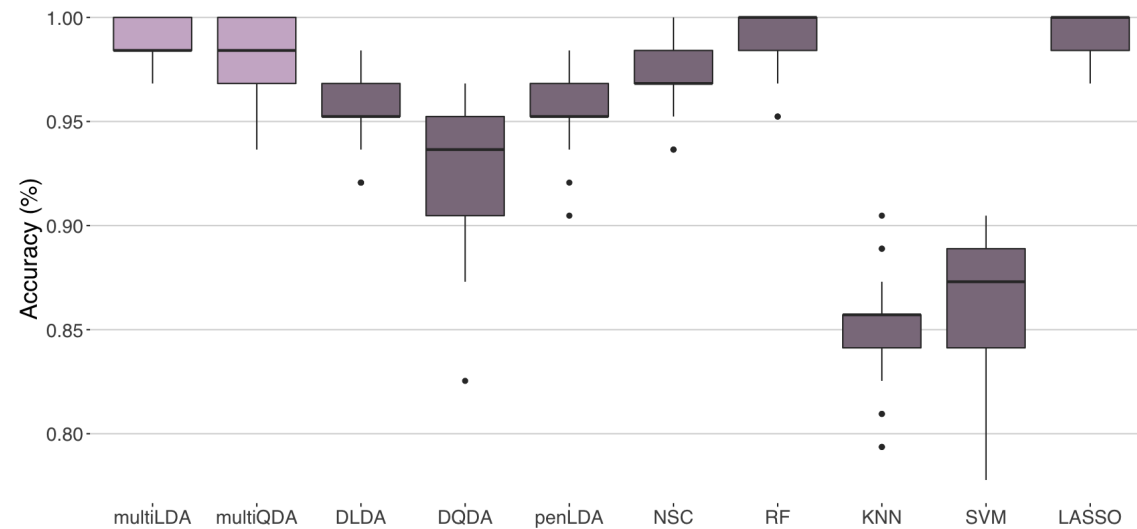
```
#remotes::install_github("sarahromanes/multiDA")  
library(multiDA)  
res <- multiDA(y = SRBCT$y,  
               X = SRBCT$X,  
               penalty = "EBIC",  
               equal.var = TRUE,  
               set.options = "exhaustive")
```

We can then examine the class groupings using the `plot()` method for `multiDA`:

```
plot(res, ranks= 1)
```



Compare performance - 100 trial, 5 fold CV





Made by a human with a computer

Slides at <https://iml.numbat.space>.

Code and data at <https://github.com/numbats/iml>.

Created using R Markdown with flair by [xaringan](#), and [kunoichi](#) (female ninja) style.



This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](#).

