

Lecture 10: Regularised regression (cont'd)

Felix Held, Mathematical Sciences

MSA220/MVE441 Statistical Learning for Big Data

6th May 2021



CHALMERS
UNIVERSITY OF TECHNOLOGY



UNIVERSITY OF GOTHENBURG

Regularisation in classification

Recall: Regularised Discriminant Analysis (RDA)

Given training samples (i_l, \mathbf{x}_l) , quadratic DA models

$$p(\mathbf{x}|i) = N(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \quad \text{and} \quad p(i) = \pi_i$$

Estimates $\hat{\boldsymbol{\mu}}_i$, $\hat{\boldsymbol{\Sigma}}_i$ and $\hat{\pi}_i$ are straight-forward to find,...

...but evaluating the normal density requires inversion of $\hat{\boldsymbol{\Sigma}}_i$. If it is (near-)singular, this can lead to **numerical instability**.

Regularisation can help here:

- ▶ Use $\hat{\boldsymbol{\Sigma}}_i = \hat{\boldsymbol{\Sigma}}_i^{\text{QDA}} + \lambda \hat{\boldsymbol{\Sigma}}^{\text{LDA}}$ for $\lambda > 0$
- ▶ Use LDA (i.e. $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}$) and $\hat{\boldsymbol{\Sigma}} = \hat{\boldsymbol{\Sigma}}^{\text{LDA}} + \lambda \boldsymbol{\Delta}$ for $\lambda > 0$ and a diagonal matrix $\boldsymbol{\Delta}$

Recall: Naive Bayes LDA

Naive Bayes LDA means that we assume that $\hat{\Sigma} = \hat{\Delta}$ for a diagonal matrix $\hat{\Delta}$. The diagonal elements are estimated as

$$\hat{\Delta}^{(j,j)} = \frac{1}{n - K} \sum_{i=1}^K \sum_{i_l=i} (\mathbf{x}_l^{(j)} - \hat{\mu}_i^{(j)})^2$$

which is the **pooled within-class variance**.

Classification is performed by predicting the class with the maximal **discriminant function** value

$$\begin{aligned} \delta_i(\mathbf{x}) &= -\frac{1}{2}(\mathbf{x} - \hat{\mu}_i)^\top \hat{\Delta}^{-1}(\mathbf{x} - \hat{\mu}_i) + \log(\hat{\pi}_i) \\ &= -\frac{1}{2} \left\| \hat{\Delta}^{-1/2}(\mathbf{x} - \hat{\mu}_i) \right\|_2^2 + \log(\hat{\pi}_i) \end{aligned}$$

where $(\hat{\Delta}^{-1/2})^{(i,i)} = 1/\sqrt{\hat{\Delta}^{(i,i)}}$.

Shrunk centroids (I)

In high-dimensional problems ($p > n$), centroids will

- ▶ contain noise
- ▶ be hard to interpret when all variables are active

As in regression, we would like to perform **variable selection** and **reduce noise**.

Recall: The class centroids solve

$$\hat{\mu}_i = \frac{1}{n_i} \sum_{l=i} \mathbf{x}_l = \arg \min_{\mathbf{v}} \frac{1}{2} \sum_{l=i} \|\mathbf{x}_l - \mathbf{v}\|_2^2$$

Idea: Can we perform variable selection through ℓ_1 -/lasso-style regularisation?
How can we account for varying variance in features and stabilise against noise?

Shrunk centroids (II)

Nearest shrunk centroids performs variable selection and stabilises centroid estimates by solving

$$\bar{\mu}_i = \arg \min_{\mathbf{v}} \frac{1}{2} \sum_{i_l=i} \|(\hat{\Delta} + s_0 \mathbf{I}_p)^{-1/2}(\mathbf{x}_l - \mathbf{v})\|_2^2 + \lambda n_i m_i \|\mathbf{v} - \hat{\mu}_T\|_1$$

where $s_0 = \text{median}(\hat{\Delta}^{(1,1)}, \dots, \hat{\Delta}^{(p,p)})$, $m_i = \sqrt{\frac{1}{n_i} - \frac{1}{n}}$ and $\hat{\mu}_T = \frac{1}{n} \sum_l \mathbf{x}_l$.

- ▶ Penalises distance of class centroid to the overall centroid μ_T
- ▶ $\hat{\Delta} + s_0 \mathbf{I}_p$ is the diagonal regularised within-class covariance matrix. Features that are highly variable across samples are scaled down (**interpretability**)
- ▶ $n_i m_i$ scales λ in case of unequal class sizes

Shrunk centroids (III)

The solution for component j can be derived as

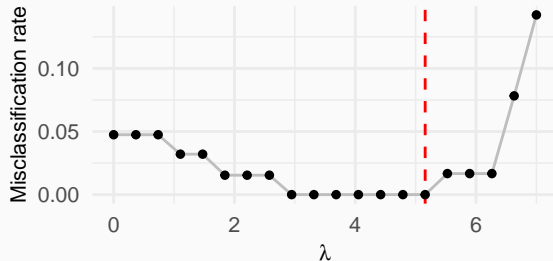
$$\bar{\mu}_i^{(j)} = \hat{\mu}_T^{(j)} + m_i(\hat{\Delta}^{(j,j)} + s_0) \text{ST}(\mathbf{t}_i^{(j)}, \lambda) \quad \text{where} \quad \mathbf{t}_i^{(j)} = \frac{\hat{\mu}_i^{(j)} - \hat{\mu}_T^{(j)}}{m_i(\hat{\Delta}^{(j,j)} + s_0)}.$$

Note: λ is a tuning parameter and has to be determined through e.g. cross-validation.

- ▶ Typically, misclassification rate improves first with increasing λ and declines for too high values
- ▶ The larger λ the more components will be equal to the respective component of the overall centroid.

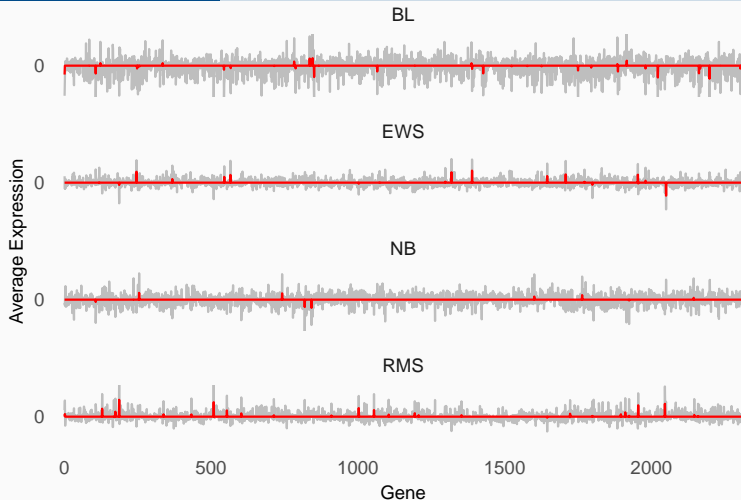
Application of nearest shrunken centroids (I)

A gene expression data set with $n = 63$ and $p = 2308$. There are four classes (cancer subtypes) with $n_{\text{BL}} = 8$, $n_{\text{EWS}} = 23$, $n_{\text{NB}} = 12$, and $n_{\text{RMS}} = 20$.



5-fold cross-validation curve and largest λ that leads to minimal misclassification rate

Application of nearest shrunken centroids (II)



Grey lines show the original centroids and red lines show the shrunken centroids

Extensions of the lasso

The lasso and groups of highly correlated variables

- ▶ The lasso does not handle groups of highly correlated variables well.
- ▶ **Example:** Two groups of highly correlated variables, e.g.

$$\mathbf{X} \sim N(\mathbf{0}, \Sigma) \quad \text{where} \quad \Sigma = \begin{pmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \Sigma_1 \end{pmatrix} \in \mathbb{R}^{200 \times 200},$$

where

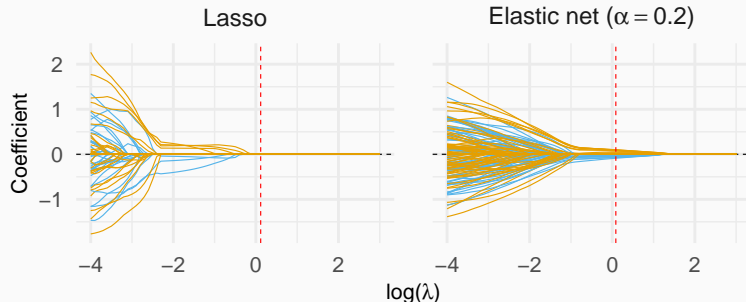
$$\Sigma_1 \in \mathbb{R}^{100 \times 100}, \quad \Sigma_1^{(i,i)} = 1.04 \quad \text{and} \quad \Sigma_1^{(i,j)} = 1, \quad i \neq j.$$

The response is generated for $n = 100$ samples as

$$\mathbf{y} = \mathbf{x}_1 - \mathbf{x}_{102} + \varepsilon \quad \text{where} \quad \varepsilon \sim N(\mathbf{0}, 4\mathbf{I}_p).$$

- ▶ **Expectation:** Since the predictors in each group are strongly correlated, all could be considered equally as predictors.
- ▶ **Possible caveat:** The lasso makes a sparsity assumption and tries to set as many coefficients to zero as possible.

The lasso and groups of highly correlated variables in practice



- ▶ At optimal λ the lasso selects 8 non-zero coefficients 0 of which were in the true coefficient vector.
 - ▶ Very precise but 'wrong' estimates.
- ▶ An alternative algorithm, the **elastic net** estimates 95 non-zero coefficients. (44 in the 1st group and 51 in the 2nd group, group-wise close coefficients)
 - ▶ 'Shares' responsibility among correlated variables

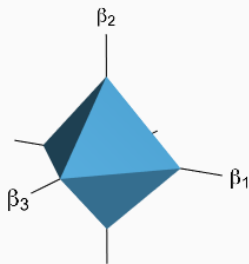
The elastic net (I)

The **elastic net** solves the problem

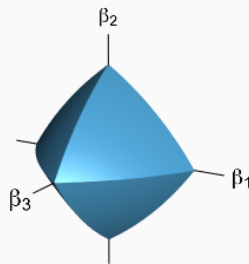
$$\arg \min_{\beta} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \left(\frac{1-\alpha}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)$$

striking a balance between lasso (**variable selection**) and ridge regression (**grouping of variables**)

Lasso



Elastic net ($\alpha = 0.7$)



- ▶ The solution can be found through **cyclic coordinate descent**
- ▶ α is an additional tuning parameter that should be determined by cross-validation
- ▶ The lasso and ridge regression are special cases of the elastic net ($\alpha = 1$ and $\alpha = 0$, respectively).

Explicitly adding groups to the lasso

- ▶ The lasso in its original formulation considers each variable separately
- ▶ Groups in data can form through e.g.
 - ▶ Correlation
 - ▶ Categorical variables in dummy encoding
 - ▶ Domain-knowledge (e.g. genes in the same signal pathway, signals that only appear in groups in a compressed sensing problem,...)
- ▶ Ideally the whole group is either present or not
- ▶ The elastic net can find groups, but only does so for highly correlated variables and without external influence. Sometimes more control is necessary.

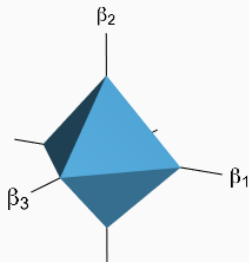
The group lasso (I)

The **group lasso** solves the problem

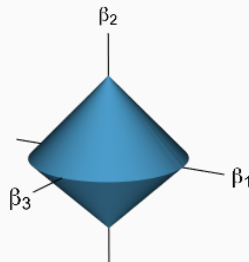
$$\arg \min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{k=1}^K \|\mathbf{B}_k\|_2$$

where \mathbf{B}_k is a vector of coefficients β_i for the k -th group. Note that $\|\beta_i\|_2 = |\beta_i|$ for singleton groups.

Lasso

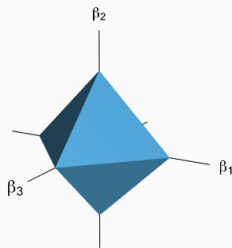


Group lasso ($\{\beta_1, \beta_3\}, \{\beta_2\}$)

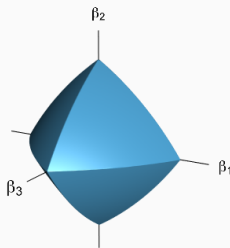


Comparison: Lasso, elastic net and group lasso

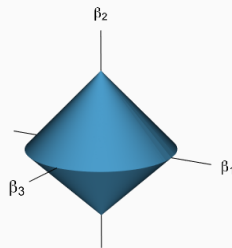
Lasso



Elastic net ($\alpha = 0.7$)



Group lasso ($\{\beta_1, \beta_3\}, \{\beta_2\}$)



- ▶ The lasso sets variables exactly to zero either on a corner or along an edge.
- ▶ The elastic net similarly sets variables exactly to zero on a corner or along an edge. The curved edges encourage remaining coefficients to be closer together.
- ▶ The group lasso has actual information about groups of variables. It encourages whole groups to be zero or non-zero with similar coefficients.

Penalisation in GLMs

Penalisation can also be used in generalised linear models (GLMs), e.g. to perform **sparse logistic regression**.

Given $p(y|\boldsymbol{\beta}, \mathbf{x})$ the log-likelihood of the model is

$$\mathcal{L}(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) = \sum_{l=1}^n \log(p(y_l|\boldsymbol{\beta}, \mathbf{x}_l))$$

Instead of penalising the minimisation of the residual sum of squares (RSS), the **minimisation of the negative log-likelihood is penalized**, i.e.

$$\arg \min_{\boldsymbol{\beta}} -\mathcal{L}(\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}) + \lambda \|\boldsymbol{\beta}\|_1$$

Note: If $p(y|\boldsymbol{\beta}, \mathbf{x})$ is Gaussian and the linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ is assumed, this is equivalent to the lasso.

Sparse logistic regression

Recall: For logistic regression with $i_l \in \{0, 1\}$ it holds that

$$p(1|\boldsymbol{\beta}, \mathbf{x}) = \frac{\exp(\mathbf{x}^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}^\top \boldsymbol{\beta})} \quad \text{and} \quad p(0|\boldsymbol{\beta}, \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{x}^\top \boldsymbol{\beta})}$$

and the penalised minimisation problem becomes

$$\arg \min_{\boldsymbol{\beta}} - \sum_{l=1}^n (i_l \mathbf{x}_l^\top \boldsymbol{\beta} - \log(1 + \exp(\mathbf{x}^\top \boldsymbol{\beta}))) + \lambda \|\boldsymbol{\beta}\|_1$$

- ▶ The minimisation problem is still convex, but non-linear in $\boldsymbol{\beta}$. Iterative quadratic approximations combined with coordinate descent can be used to solve this problem.
- ▶ Another way to perform sparse classification (like e.g. nearest shrunken centroids)

Sparse multi-class logistic regression

In multi-class logistic regression with $i_l \in \{1, \dots, K\}$, there is a matrix of coefficients $\mathbf{B} \in \mathbb{R}^{p \times (K-1)}$ and it holds for $i = 1, \dots, K-1$ that

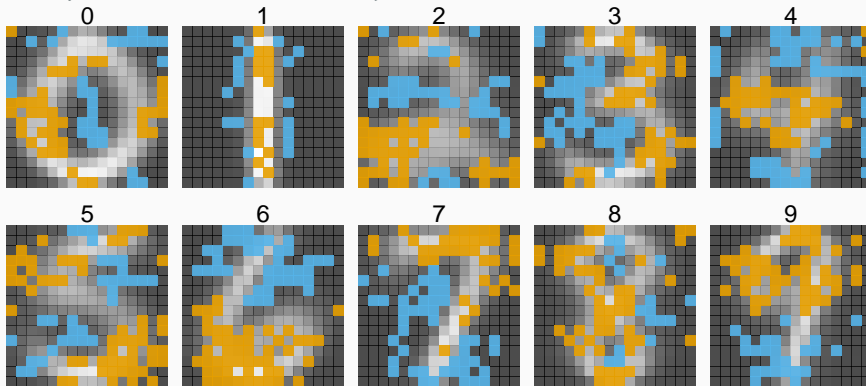
$$p(i|\mathbf{B}, \mathbf{x}) = \frac{\exp(\mathbf{x}^\top \boldsymbol{\beta}_i)}{1 + \sum_{j=1}^{K-1} \exp(\mathbf{x}^\top \boldsymbol{\beta}_j)} \quad \text{and} \quad p(K|\mathbf{B}, \mathbf{x}) = \frac{1}{1 + \sum_{j=1}^{K-1} \exp(\mathbf{x}^\top \boldsymbol{\beta}_j)}$$

- ▶ As in two-class case, the absolute value of each entry in \mathbf{B} can be penalised.
- ▶ Another possibility is to use the group lasso on all coefficients for one variable, i.e. penalise with $\|\mathbf{B}_{j\cdot}\|_2$ for $j = 1, \dots, p$.

Example for sparse multi-class logistic regression

MNIST-derived zip code digits ($n = 7291$, $p = 256$)

Sparse multi-class logistic regression was applied to the whole data set and the penalisation parameter was selected by 10-fold CV.



Orange tiles show positive coefficients and blue tiles show negative coefficients. Class averages are shown in the background.

Take-home message

- ▶ Penalisation methods are not only restricted to regression, also applicable to classification
- ▶ Sparsity is a very important concept when interpretability of models is important
- ▶ Many extensions to the lasso exist, which make it more suitable for a variety of different situations