

## Exercise 6

**Deadline: 28.01.2022, 16:00.**

Ask questions to [#ask-your-tutor-christoph](#)

This exercise is dedicated to regularized regression.

### Regulations

Please hand in your solution as a Jupyter notebook `regularized-regression.ipynb`, accompanied with exported `regularized-regression.html`. Zip all files into a single archive `ex06.zip` and upload this file to your assigned tutor on MaMPF before the given deadline.

**Note:** Each team creates only a single upload, and all team members must *join* it as described in the MaMPF documentation at <https://mampf.blog/zettelabgaben-fur-studierende/>.

**Important:** Make sure that your MaMPF name is the same as your name on Muesli. We now identify submissions purely from the MaMPF name. If we are unable to identify your submission you will not receive points for the exercise!

### 1 Bias and variance of ridge regression (8 points)

Ridge regression solves the regularized least squares problem

$$\hat{\beta}_\tau = \operatorname{argmin}_\beta (y - X\beta)^\top (y - X\beta) + \tau \beta^\top \beta$$

with regularization parameter  $\tau \geq 0$ . Regularization introduces some bias into the solution in order to achieve a potentially large gain in variance. Assume that the true model is  $y = X\beta^* + \epsilon$  with zero-mean Gaussian noise  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  and centered features  $\frac{1}{N} \sum_i X_i = 0$  (note that these assumptions imply that  $y$  is also centered in expectation). Prove that expectation and covariance matrix of the regularized solution (both taken over all possible training sets of size  $N$ ) are then given by

$$\mathbb{E}[\hat{\beta}_\tau] = S_\tau^{-1} S \beta^* \quad \operatorname{Cov}[\hat{\beta}_\tau] = S_\tau^{-1} S S_\tau^{-1} \sigma^2$$

where  $S$  and  $S_\tau$  are the ordinary and regularized scatter matrices:

$$S = X^\top X \quad S_\tau = X^\top X + \tau \mathbb{I}_D$$

Notice that expectation and covariance reduce to the corresponding expressions of ordinary least squares (as derived in the lecture) when  $\tau = 0$ :

$$\mathbb{E}[\hat{\beta}_{\tau=0}] = \beta^* \quad \operatorname{Cov}[\hat{\beta}_{\tau=0}] = S^{-1} \sigma^2$$

Since  $S_\tau$  is greater than  $S$  (in any norm), regularization has a shrinking effect on both expectation and covariance.

### 2 Denoising of a CT image (11 points)

In the last task of exercise 5, we diminished the radiation dose for our patient by using fewer projection angles. However, this introduced noise in the resulting tomogram. Ridge regression offers a simple possibility to reduce the noise level<sup>1</sup>.

<sup>1</sup>You will notice that the improvements don't look too impressive. Humans tend to prefer the noisy, but unbiased version of the image, because the eye's built-in noise suppression is way more powerful than ridge regression. In other words, the eye and the algorithm have different bias-variance trade-offs.

To apply regularization to tomographic reconstruction, the  $D \times D$  diagonal matrix  $\sqrt{\tau} \mathbb{I}_D$  must be appended at the bottom of the weight matrix  $X$  from exercise 5:

$$X' = \begin{bmatrix} X \\ \sqrt{\tau} \mathbb{I}_D \end{bmatrix}$$

Correspondingly, vector  $y$  must be extended with  $D$  zeros.

Add a new parameter to the function `construct_X()` to activate regularization

`X = construct_X(M, alphas, Np = None, tau=0)`

For `tau=0`, the function shall behave exactly like the original version from exercise 5. If `tau>0`, it shall return the augmented matrix  $X'$ . Reconstruct the tomogram for 64 angles with  $\tau = 0, 1, 10, 100, 1000, 10000$  and display the resulting CT images. Find the value of  $\tau$  with the best compromise of image sharpness and noise.

Compare the ridge regression results with Gaussian filtering, another popular noise reduction technique, which replaces each pixel with a weighted average of its neighbors (see [https://en.wikipedia.org/wiki/Gaussian\\_filter](https://en.wikipedia.org/wiki/Gaussian_filter) for more details). A suitable implementation is provided by function `gaussian_filter` in module `scipy.ndimage.filters`. Apply `gaussian_filter` with filter sizes `sigma = 1, 2, 3, 5, 7` to the CT image with  $\tau = 0$ , display the results and comment on the similarities and differences between ridge regression and Gaussian filtering.

### 3 Automatic feature selection for regression

Sparse regression regularizes the solution such that unimportant elements of  $\hat{\beta}$  are set to zero. The corresponding columns of  $X$  could be dropped from the matrix without any effect on the solution – the remaining columns are obviously more useful as an explanation for the response  $y$ . Therefore, sparse regression can be interpreted as a method for *relevant* feature identification.

In exercise 2, you implemented dimension reduction from a full image of a digit to just two features by selecting important pixels manually<sup>2</sup>. In this exercise, we will automatically find relevant pixels by means of sparse regression.

#### 3.1 Implement Orthogonal Matching Pursuit (5 points)

“Orthogonal Matching Pursuit”<sup>3</sup> is a simple greedy sparse regression algorithm. It approximates the exact algorithms for least squares with  $L_0$  or  $L_1$  regularization (cf. the lecture for details) and is defined as:

**Initialization:**

- Inputs:  $X \in \mathbb{R}^{N \times D}$ ,  $y \in \mathbb{R}^N$ ,  $T > 0$  (the desired number of non-zero elements in the final solution  $\hat{\beta}^{(T)}$ )
- Define the initial sets of active resp. inactive columns as  $A^{(0)} = \emptyset$  and  $B^{(0)} = \{j : j = 1 \dots D\}$ .
- Define the initial residual  $r^{(0)} = y$ .

**Iteration:** for  $t = 1 \dots T$  do

1. Find the inactive column that has maximal correlation with the current residual:

$$j^{(t)} = \operatorname{argmax}_{j \in B} \left| X_j^\top \cdot r^{(t-1)} \right|$$

<sup>2</sup>or computing your own features

<sup>3</sup>not to be confused with its simpler cousin “Matching Pursuit”

2. Move  $j^{(t)}$  from the inactive set  $B$  to the active set  $A$ .
3. Form the active matrix  $X^{(t)}$  consisting of all currently active columns of  $X$ .
4. Solve the least squares problem

$$\hat{\beta}^{(t)} = \operatorname{argmin}_{\beta} (y - X^{(t)}\beta)^{\top} (y - X^{(t)}\beta)$$

5. Update the residual  $r^{(t)} = y - X^{(t)}\hat{\beta}^{(t)}$ .

Implement this algorithm as a function

```
solutions = omp_regression(X, y, T)
```

which returns  $\hat{\beta}^{(t)}$  for  $t = 1 \dots T$  as a  $D \times T$  matrix, i.e. the inactive elements in each solution are not dropped, but explicitly set to zero.

### 3.2 Classification with sparse LDA (8 points)

We again use the `digits` dataset of scikit-learn and classify ‘1’ vs. ‘7’. For balanced training sets, LDA can be formulated as a least squares problem, where the rows  $X_i$  are the flattened images of each digit, and

$$y_i = \begin{cases} 1 & \text{if instance } i \text{ is a ‘1’} \\ -1 & \text{if instance } i \text{ is a ‘7’} \end{cases}$$

are the desired responses. Now execute `omp_regression()` with sufficiently big  $T$  to get the sequence of sparse LDA solutions for  $t = 1 \dots T$ .

Report the error rate on the test set for  $t = 1 \dots T$ . How many pixels should be used for acceptable error rates? Is it necessary/beneficial to standardize the data before training and testing?

Visualize in which order the pixels are switched to *active* as  $t$  increases, and show if a pixel votes in favour or against class ‘1’. What is a good criterion for this distinction? Compare these results with your hand-crafted feature selection in exercise 2 – did you select the same pixels?

### 3.3 One-against-the-rest classification (8 points)

A *one-against-the-rest classifier* is actually a collection of  $C$  classifiers, one for each class  $k$ . To construct it, one uses the given training data to create  $C$  auxiliary training sets with target responses

$$y_i^{(k)} = \begin{cases} 1 & \text{if instance } i \text{ has class } k \\ -1 & \text{if instance } i \text{ has class } \neq k \end{cases}$$

The positive subset contains all instances of class  $k$ , and the negative subset an equally sized random selection from all other classes. Use these auxiliary training sets to train a sparse LDA for the digits  $k = 0 \dots 9$  with suitable  $T$  (all classifiers should use the same  $T$ ).

At testing, a new instance is rated by all  $C$  classifiers and assigned to the class whose classifier returns the largest response, or to the additional class “unknown” if all classifiers have negative response. Report the confusion matrix of your predictor on the test set.

Why is it useful to introduce class “unknown”? Parameterize your classifier such that a few test instances are actually assigned to this class and visualize these instances. What do you observe?