

1. Boston Housing

Statisticians collected data from the Boston area to search for a **formula** giving the – approximated – price of an estate depending on attributes like: *crime rate (CRIM)*, *prop. of resid. land (ZN)*, *prop. of business, etc.*¹ Due to problems met when assessing some variables – including the price of the house – the data-set is quite noisy.

Your task is to build a model that is successful in predicting the house price. You should use a generalised linear model with the following setup:

$$\mathbf{x}_n = [x_{n1}, x_{n2}, \dots, x_{nd}], \quad y_n = \text{declared price}, \quad \mathbf{y} = [y_1, \dots, y_N]^T$$

$$f(x_n|\boldsymbol{\theta}) = \sum_{k=1}^K \theta_k \phi_k(\mathbf{x}_n) = \boldsymbol{\phi}_n \boldsymbol{\theta} \quad \text{with } \phi_k \text{ a function base of your choice}$$

with the notations: $\boldsymbol{\theta} \stackrel{\text{def}}{=} [\theta_1, \dots, \theta_K]^T, \quad \boldsymbol{\phi}_n = [\phi_1(\mathbf{x}_n), \dots, \phi_K(\mathbf{x}_n)]$

Consider using the M.A.P. procedure where the solution is:

$$\boldsymbol{\theta}^* = \left(\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \frac{1}{\sigma_0^2} \mathbf{I}_K \right)^{-1} \boldsymbol{\Phi}^T \mathbf{y}$$

To assess the **error of the model**, use the *leave-one-out procedure*, i.e. design a testing procedure that measures the error by:

- Using a data set with *a single missing element* (\mathbf{x}_n, y_n) , denote it with \mathcal{D}^{-n} ;
- Computing the best model $\boldsymbol{\theta}_{-n}^*$ based on \mathcal{D}^{-n} ;
- Compute the square error for the data left-out: $LOO(n) = (y_n - f(\mathbf{x}_n|\boldsymbol{\theta}_{-n}^*))^2$.

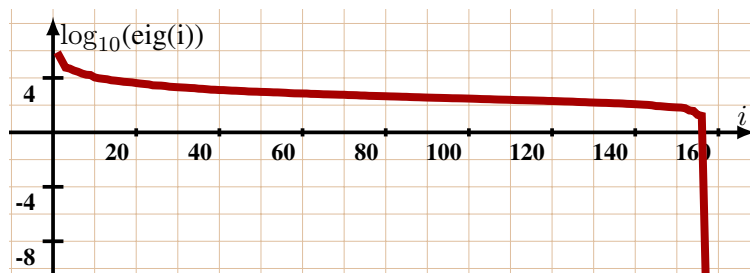
and finally compute $LOO(\mathcal{D}) = \sum_n LOO(n)$

Analyse the results of the experiments: display the errors as a function of the *number of basis functions* and the *a-priori distribution of the parameters*, σ_0^2 .

2. Yale faces²

The Yale face dataset contains 165 images of subjects in different circumstances. The images are of size [243, 320], thus the number of attributes for each observation is 77760 and there are altogether 165 images.

The plotted eigenvectors – on the right – indicate that there are close to 160 basis elements needed to *compactly* represent the data-set.



- Implement procedures to reduce this number by modifying the images with a minimum “loss” of information and such that the eigenvectors will drop much quicker – making possible the implementation of a simple dimensionality reduction technique. (hint: try to remove – algorithmically – “empty” zones from the top, bottom, left, right. Since the regions are shaded, you should detect first the position of the e.g. nose-eyes)
- Label the data as **those-who-wear-glasses** and **those-who-not** and use the leave-one-out (LOO) procedure defined before to test whether the **original data-set** or the **reduced one** produces better results.

You can access the Yale data-set at http://www.cs.ubbcluj.ro/~csatol/prob_datamin/data/yale.mat and additional M-files are available (like visualising `vis(yale, ind, [243, 320])`) in the http://www.cs.ubbcluj.ro/~csatol/prob_datamin/matlab directory.

¹ See <http://lib.stat.cmu.edu/S/Harrell/data/descriptions/boston.html> for the data and the full description of the variables. The data itself is e.g. at <http://lib.stat.cmu.edu/datasets/boston>.

² Available from: <http://cvc.yale.edu/projects/yalefaces/yalefaces.html>